

MULTIOBJECTIVE TIME SERIES MATCHING AND CLASSIFICATION

PHILIPPE ESLING

Philosophical Doctor (PhD) in Acoustics, Signal Processing and
Computer Science

Equipe représentations musicales

Institut de Recherche et Coordination Acoustique Musique (IRCAM)

Université Pierre et Marie Curie

March 2012



Ohana means family.
Family means nobody gets left behind, or forgotten.

— Lilo & Stitch

Dedicated to the loving memory of Rudolf Miede.

1939 – 2005

ABSTRACT

Millions of years of genetic evolution have shaped our auditory system, raising our way of listening to a form of art. Despite a somehow limited frequency spectrum, we are able to achieve an excellent and flexible discrimination of events. These unique capacities originate from the ability of our brain to organize our perception of sounds and music. We can process several conflicting scales simultaneously, thus constructing a multidimensional structure of perception. Furthermore, even if time is an ubiquitous and complex concept, humans have a natural capacity to extract meaningful knowledge from the shape of temporal structures. The onset of this study was, therefore, to explore these temporal and perceptual aspects in order to create a framework for generating musical orchestrations.

We show that by drawing inspiration from our musical perception and gaining insights from these mechanisms to drive our choices of algorithms, we can create innovative and powerful approaches for generic querying and classification, far outside the realm of musical problematics. First, by trying to emulate our multiobjective perception of temporal structures, we propose a framework called *MultiObjective Time Series* (MOTS) matching. We formally state this novel problem and provide an efficient algorithm to solve it. Based on this approach, we are able to introduce two innovative audio querying paradigms. We then examine their effectiveness and usability through extensive user studies. We prove the validity of our proposal by studying the perception of the temporal evolution of conflicting higher-level audio features. We reveal the concept of multidimensional *directions of listening* that forms in the brain. We show that these directions are consistent through various tasks and unique to each person. We further propose a novel and flexible classification model based on the *hypervolumes dominated* by different classes, called *HyperVolume-MOTS* (HV-MOTS) classification. Instead of trying to consider the position of an element with respect to the various classes, this framework studies the behavior of each class with respect to the input through the distribution and spread over the optimization space. We show that the multiobjective flexibility inspired by our musical perception produces a classification paradigm that outperforms state-of-the-art methods on a wide range of scientific problems such as EEG analysis, climatology, medical diagnosis, character recognition and robotics. We present a comparison of this classification paradigm to traditional classifiers such as Nearest-Neighbor, Nearest-Center or Support Vector Machines. We then perform a comprehensive and thorough evaluation of our new approach and demonstrate its superiority on a wide range of datasets. We also show several applications of this scheme and study its weaknesses and strengths in each case. We present our main finding in which this method allows to construct a biometric identification system based on the sounds produced by heartbeats. We specifically develop for this problem a novel set of features based on the Stockwell transform and inspired by research in musical analysis. We show that we can accurately identify human beings through the sounds their heart produce. Our system obtains error rates equivalent to other biometrics such as face or speech recognition. These findings are supported by the largest heart sounds dataset ever collected, including the Mars500 isolation study.

Finally, we show how all this knowledge gained allows to come back to our initial artistic problematics of musical orchestration. We consider the problem of generating

orchestral sound mixtures that can closely approximate any given audio signal. While performing this reconstruction, we avoid mixing the similarities into a single measure, but rather use an advanced search algorithm based on the MOTS framework, called optimal warping. This allows us to obtain a set of efficient solutions that provide various compromises among spectral objectives. This algorithm performs a morphological segmentation procedure based on the variation of entropy. We then present several musical applications and interfaces that result from our various generic findings.

RÉSUMÉ

Plusieurs millions d'années d'évolution génétique ont façonné notre système auditif, élevant ainsi notre écoute au rang d'un art. Malgré un spectre de fréquences perçues quelque peu limité, nous sommes en mesure d'effectuer une discrimination précise et flexible des événements auditifs. Ces capacités uniques proviennent de la capacité qu'a notre cerveau à organiser notre perception des sons et de la musique. Nous pouvons ainsi traiter simultanément plusieurs échelles de perception contradictoires, par la construction d'une structure multidimensionnelle de la perception. De plus, même si le temps est un concept omniprésent et complexe, les êtres humains ont une capacité inhérente à extraire une structure cohérente à partir de formes temporelles. Le point de départ de notre travail était donc d'étudier ces aspects temporels et perceptuels pour la création d'un système de génération d'orchestration musicale.

Nous montrons qu'en s'inspirant de cette perception musicale et en émulant ces mécanismes dans nos choix algorithmiques, nous sommes en mesure de créer des approches novatrices et efficaces de recherche et de classification générique, dépassant largement le cadre des problématiques musicales. Tout d'abord, en essayant d'imiter le caractère multi-objectif de notre perception des structures temporelles, nous proposons un cadre de recherche appelé *MultiObjective Time Series* (MOTS). Nous commençons par définir formellement ce nouveau problème et proposons un algorithme efficace pour le résoudre. Sur la base de cette approche, nous sommes en mesure d'introduire deux paradigmes innovants de recherche sur les fichiers audio. Nous étudions l'efficacité et la facilité d'utilisation de ces paradigmes grâce à des études utilisateurs à grande échelle. Grâce à cette étude, nous prouvons également la validité de notre proposition en analysant la perception d'évolutions temporelles conflictuelles sur des descripteurs audio de haut niveau. Nous exposons ainsi le concept de *directions d'écoute* multidimensionnelles qui prends naissance dans notre perception. Nous montrons que ces directions sont consistantes à travers plusieurs tâches mais également uniques à chaque personne. Après cette validation, nous introduisons un nouveau paradigme flexible de classification basé sur les *hypervolumes dominés* par les différentes classes, appelé *HyperVolume-MOTS* (HV-MOTS). Contrairement aux paradigmes classiques qui étudient la position d'un élément par rapport aux différentes classes existantes, notre système étudie le comportement de la classe entière par rapport à l'élément à travers la distribution et la diffusion d'une classe sur l'espace d'optimisation. Nous montrons que la flexibilité multi-objective inspirée par notre perception musicale produit un paradigme de classification qui surpasse les méthodes de l'état de l'art sur un large éventail de problèmes scientifiques tels que l'analyse EEG, la climatologie, le diagnostic médical, la reconnaissance de caractères et la robotique. Nous fournissons une comparaison de ce paradigme par rapport aux classificateurs classiques tels que le Nearest-Neighbor, Nearest-Center ou Support Vector Machines. Nous effectuons ensuite une évaluation exhaustive et approfondie de notre nouvelle approche et démontrons sa supériorité sur un large ensemble de données. Nous montrons en outre plusieurs applications permettant d'étudier de manière plus détaillée les forces et faiblesses de notre proposition. Nous présentons l'application principale de cette méthode dans laquelle elle permet de construire un système d'identification biométrique basée sur les sons produit par les battements de coeur. En particulier, nous développons pour ce problème un nouvel ensemble de descripteurs basés sur la transformée de

Stockwell et inspiré par la recherche en analyse musicale. Nous montrons que nous pouvons identifier avec précision les êtres humains à travers les sons que produit leur cœur et que nous atteignons des taux d'erreur équivalents à d'autres caractéristiques biométriques telles que la reconnaissance vocale. Ces résultats sont confirmés par le plus grand ensemble de données de sons cardiaques jamais recueillies, comprenant également l'étude d'isolation Mars500 effectuée par l'Agence Spatiale Européenne.

Enfin, nous montrons comment toute cette connaissance acquise permet de revenir à nos problématiques artistiques originales d'orchestration musicale. Nous étudions ainsi le problème de la génération de mélanges sonores orchestraux imitant au mieux une cible audio donnée. En effectuant cette reconstruction, nous évitons de mélanger la similarité en une mesure de distance unique et nous utilisons un nouvel algorithme de recherche basé sur le cadre MOTS appelé *Optimal Warping*. Cette approche nous permet ainsi d'obtenir un ensemble de solutions efficaces qui offrent différents compromis entre les objectifs spectraux. Cet algorithme effectue une segmentation morphologique basée sur l'analyse de la variation d'entropie des séries temporelles. Nous présentons enfin plusieurs interfaces et applications musicales qui résultent de nos travaux.

PUBLICATIONS

Parts of this thesis along with some ideas and figures have appeared previously in the following journal publications :

INTERNATIONAL JOURNALS

- 2012 **Esling Philippe**, Agon Carlos "Time series data mining and analysis", *ACM Computing Surveys*, vol. 46, no. 1, 2013.
- 2012 **Esling Philippe**, Agon Carlos "Multiobjective time series matching for audio classification and retrieval", *IEEE Transactions on Speech Audio and Language Processing* 2013 (Accepted - Major changes).
- 2012 Hacbarth Benjamin, Schnell Norbert, **Esling Philippe**, Schwarz Diemo "Composing Morphology: Concatenative Synthesis as an Intuitive Medium for Prescribing Sound in Time", *Contemporary Music Review* (to appear)
- 2011 Lecroq Béatrice, Lejzerowicz Franck, **Esling Philippe**, Baerlocher Loic, Farinelli Laurent, Pawlowski Jan "Ultra-deep sequencing of foraminiferal microbarcodes unveils hidden richness of early monothalamous lineages in deep-sea sediments", *Publication of the National Academy of Science*, vol.108, no.32, pp 13177-13182, August 2011.
- **Esling Philippe**, Agon Carlos "Let your heart beat and I will tell you who you are", 2013

BOOK CHAPTERS

- 2010 **Esling Philippe**, Carpentier Grégoire, Agon Carlos "Dynamic Musical Orchestration using Genetic Algorithms and a Spectro-Temporal Description of Musical Instruments", *Lecture Notes in Computer Science*, vol. 6025, *EvoApplications Part II*, 2010.

INTERNATIONAL CONFERENCES (WITH REVIEW COMITEE)

- 2010 **Esling Philippe**, Carpentier Grégoire, Agon Carlos "Dynamic Musical Orchestration using Genetic Algorithms and a Spectro-Temporal Description of Musical Instruments", *Lecture Notes in Computer Science*, vol. 6025, *EvoApplications Part II*, 2010.
- 2010 **Esling Philippe**, Agon Carlos "Composition of Sound Mixtures with Spectral Maquettes", *Proceedings of the International Computer Music Conference*, New York, USA, 2010.
- 2010 **Esling Philippe**, Agon Carlos "Composer les mélanges sonores avec les maquettes spectrales", *Actes des 10emes Journées d'Informatique Musicale*, pp. 5-15, Rennes, France, 2010

NATIONAL CONFERENCES

- 2010 **Esling Philippe**, Agon Carlos "Time series analysis, sound mixtures and orchestration", Presentation in CNRS Japanese-French Laboratory of Informatics, Tokyo University, 2010.
- 2009 **Esling Philippe**, Agon Carlos "Orchestration and Sound Mixtures", Journées Jeunes Chercheurs en Acoustique Audition et Signal, Marseille, 2009.

*When you want the food, you take the food.
When you want the car, you take the car.
When you want the woman ... you take the woman*

— El Jefe

ACKNOWLEDGMENTS

Many thanks to everybody, I love you all ! Carlos “Jefe” Agon, Karim “El Guapo” Haddad, ¹

¹ El Guapo

CONTENTS

I INTRODUCTION	1
1 The artistic problematic	2
1.1 From the past empiricism ...	2
1.2 ... through musical writing ...	3
1.3 ... towards a modern treatise ?	4
1.4 Elementary units of orchestral study	5
2 Studying elements	6
2.1 Musical atoms	6
2.2 Signal and symbolism	7
2.3 Multiplicities of timbre	7
2.3.1 Timbre and acousticians	8
2.3.2 Timbre and musicians	8
2.3.3 Timbre and computer scientists	9
2.4 Time scales continuum	9
3 Perceiving elements	11
3.1 The doors of perception	11
3.2 Complex and multifaceted similarities	11
3.2.1 Assessing visual similarity	11
3.2.2 Assessing sound similarity	13
3.3 Multidimensionality of timbre perception	14
4 Putting the pieces together	16
4.1 Rationale of this study	16
4.2 Epistemological considerations	16
4.3 Scientific contributions	16
4.4 Structure of this document	18
II ELEMENTS OF TIME AND PERCEPTION	21
5 Time series data mining	22
5.1 Introduction	22
5.2 Definitions	23
5.3 Tasks in time series data mining	24
5.3.1 Query by content	24
5.3.2 Clustering	26
5.3.3 Classification	28
5.3.4 Segmentation	29
5.3.5 Prediction	30
5.3.6 Anomaly detection	31
5.3.7 Motif discovery	32
5.4 Implementation components	33
5.4.1 Preprocessing	34
5.4.2 Representation	34
5.4.3 Similarity measure	38
5.4.4 Indexing	43
5.5 Research trends and issues	46
6 Multiobjective optimization	49

6.1	Definitions	49
6.1.1	Pareto dominance	49
6.1.2	Chebyshev norms	51
6.2	Algorithms classification	52
6.3	Applications	53
III MULTIOBJECTIVE TIME SERIES (MOTS) MATCHING		55
7	MOTS Framework	56
7.1	Problem definition	56
7.2	Comparison to multivariate matching	59
7.3	Algorithms	59
7.3.1	Multiobjective early abandon	60
7.3.2	Hyperplane search	63
7.4	Efficiency on massive databases	66
7.4.1	Comparing algorithms	66
7.4.2	Comparing datasets	69
7.5	Innovative audio querying	70
7.5.1	Content-based audio retrieval	70
7.5.2	Going beyond traditional query paradigms	71
7.6	Database structure	73
7.6.1	QBE results and representation	74
7.6.2	MultiObjective Spectral Evolution Query (MOSEQ)	75
7.6.3	Query by Vocal Imitation (QVI)	77
8	Validating the MOTS Framework	79
8.1	Audio features	80
8.2	Hypotheses	82
8.3	Protocol	84
8.3.1	Tasks	85
8.3.2	Participants	90
8.3.3	Datasets	91
8.3.4	Equipment	92
8.4	Results	92
8.4.1	Multidimensional directions of listening	93
8.4.2	Abstract multi-dimensional similarity	106
8.4.3	Usability evaluation of audio querying paradigms	108
8.4.4	Vocal control of spectral features	114
8.4.5	Impact of skills	116
8.5	Generalization	116
9	Conclusions of this part	118
IV HYPERVOLUME CLASSIFICATION (HV-MOTS)		119
10	HV-MOTS classification	120
10.1	Multi-objective classification	120
10.2	Distance-based classifiers	122
10.3	Comparison to other classifiers	123
10.4	Discussion	126
10.4.1	Past results	126
10.4.2	Advantages and drawbacks	126
10.4.3	Comparison	127
11	Large scale study	129

11.1	Datasets summary	129
11.2	Methodology	131
11.2.1	Evaluation framework	131
11.2.2	Hardware	132
11.2.3	Algorithms implementation	132
11.2.4	Reproducibility of experiments	133
11.3	Results and analysis	133
11.3.1	Dataset-wise results	133
11.3.2	Global scale analysis	134
11.4	Comparison to state-of-art results	140
11.5	Extended analysis	142
11.5.1	Selected features	142
11.5.2	Warping or resampling	142
11.5.3	The power of time	142
12	Unicity of heart sounds	146
12.1	Biometric systems	146
12.2	Heart physiology	147
12.3	Listening to the heart	149
12.3.1	<i>Where</i> to listen (pre-processing)	149
12.3.2	<i>When</i> to listen (segmentation)	150
12.3.3	<i>What</i> to listen (S-Features)	151
12.3.4	How to listen (HV-MOTS Scoring)	152
12.4	Experiments	153
12.4.1	Datasets	153
12.4.2	Evaluation methodology	155
12.4.3	Results	155
12.4.4	Comparison to existing biometrics	163
12.5	Discussion	163
13	Audio applications	166
13.1	Intelligent sound sample database	166
13.1.1	Content-based audio retrieval	166
13.1.2	Datasets	167
13.1.3	Results analysis	169
13.1.4	Comparison to state of the art	173
13.1.5	Robustness analysis	173
13.2	Sound morphology	174
13.2.1	Onset of the study	174
13.2.2	Results	176
14	Conclusions	179
V	GOING BACK TO MUSIC	181
15	Orchestration	182
15.1	On the complexity of orchestration	182
15.1.1	Combinatorial complexity	182
15.1.2	Temporal complexity	183
15.1.3	Orchestral timbre	183
15.2	How to reify orchestration	184
15.2.1	Existing systems	184
15.2.2	Discussion	185

15.3	Going further in computer-aided orchestration	186
15.3.1	Algorithmic choices	186
15.4	Abstract Temporal Orchestration - Modular Structure (ATO-MS)	187
15.4.1	Entropic segmentation procedure	188
15.4.2	Optimal warping	188
15.4.3	Comparison with Orchidee	190
15.4.4	Modular structure	192
15.4.5	Database	193
15.4.6	Interface	194
16	Other artistic applications	196
16.1	MOSEQ Interface	196
16.2	QVI Interface	196
16.3	iPad Interface	196
16.4	Spectral Maquettes	196
16.4.1	Motivation	196
16.4.2	Implementation in OpenMusic	197
17	Conclusions of this part	200
VI	CONCLUSIONS	201
18	Future work	202
18.1	The MOTS paradigm	202
18.1.1	Wider applications	202
18.1.2	Hybrid analysis	203
18.1.3	Interaction and representation	203
18.2	MOSEQ / QVI	203
18.2.1	Applications in audio workflow	203
18.2.2	Relevance feedback	204
18.3	HV-MOTS Classification	204
18.3.1	Audio applications	204
18.3.2	Scope of application	204
18.3.3	Multiobjective subsequence classification	205
18.4	Heart sounds biometry	205
18.4.1	Segmentation procedure	205
18.4.2	Features computation	205
18.4.3	Factors of influence	205
18.5	On heart diseases detection	206
18.6	Orchestration	206
18.6.1	Signal and symbolism	206
18.6.2	Intelligent music notation	207
18.6.3	Emergence phenomenon	207
18.7	Closing the gap between all worlds	208
18.7.1	Constraint inference system	208
18.7.2	Several views on constraints	210
18.8	On a time scales continuum	210
18.8.1	From micro to macro	211
18.8.2	Macro-temporal articulations	211
19	Conclusions	212
VII	APPENDIX	215
A	APPENDIX	216

A.1	HV-MOTS Datasets description	216
A.2	Unicity of heart sounds	236
A.2.1	Cardiac auscultation	236
A.2.2	S-Features	238
A.2.3	Statistical moments	239
A.2.4	Energy distribution	241
A.2.5	Peaks distribution	242
A.3	Extended hearts biometry analysis	243
A.3.1	Pre-processing	243
A.3.2	Segmentation	243
A.3.3	Beat Selection	246
A.3.4	Time series comparison	246
A.3.5	Decision influence	250
BIBLIOGRAPHY		255

LIST OF FIGURES

Figure 1	Different levels of complexity appear in the study of music, where the processing of musical pieces is performed on several scales simultaneously. The various potential combinations of instruments induce a <i>combinatorial complexity</i> . The <i>macro-temporal evolution</i> of these musical structures strongly condition our perception. However, if we focus on a single musical atom, it also embeds <i>micro-temporal</i> evolutions of its spectral properties that are perceived in a <i>multidimensional</i> fashion. 6
Figure 2	The doors of our visual perception. We are able to perceive only a very small fraction of the electromagnetic spectrum. However, even in this narrow perceivable part we can still differentiate between millions of colors and perceive extremely complex and detailed spatial structures. 12
Figure 3	The doors of our auditory perception. As for visual perception, we can only perceive a very small fraction of the acoustic spectrum. However, in this narrow portion of the air vibrations, we are able to differentiate hundreds of pitches and perceive complex temporal evolutions and structures. Finally, even if we take a fixed portion of this acoustic spectrum (supposedly instruments at the same fixed pitch), we are still able to differentiate between all these sources based on their various properties. 12
Figure 4	The assessment of <i>visual</i> similarity in a set of elementary units. Given this set of elements, we are faced between evaluating their similarity based on the <i>shapes</i> of elements such as grouping G_1 or rather based on their <i>colors</i> such as grouping G_2 . 13
Figure 5	The assessment of sound similarity given its most elementary unit, a synthesized sinusoidal signal. For this signal, we can define the temporal function of its <i>amplitude</i> ("loudness") and its <i>frequency</i> ("pitch"). Therefore, we define two very simple temporal functions for each feature. The loudness $\mathcal{A}(t)$ can either be set to \mathcal{A}^1 or \mathcal{A}^2 and the pitch $\mathcal{F}_0(t)$ can either be set to \mathcal{F}_0^1 or \mathcal{F}_0^2 . Given these two possibilities for each feature, we can easily synthesize the set $\mathcal{S}_{\mathcal{A}}^{\mathcal{F}_0}$ of four different sounds $\mathcal{S}_1^1, \mathcal{S}_2^1, \mathcal{S}_1^2$ and \mathcal{S}_2^2 . If we try to assess the similarity between elements inside this set, it seems unfeasible to choose if we should to group elements based on the evolution of their <i>loudness</i> properties such as grouping G_1 or rather based on the temporal evolution of their <i>pitch</i> such as grouping G_2 . 15
Figure 6	The <i>epistemological loop</i> of this study. 17

- Figure 7 Diagram of a typical query by content task represented in a 2-dimensional search space. Each point in this space represents a series whose coordinates are associated with its features. (a) When a query is entered into the system, it is first transformed into the same representation as that used for other datapoints. Two types of query can then be computed. (b) A ϵ -range query will return the set of series that are within distance ϵ of the query. (c) A K-Nearest Neighbors query will return the K points closest to the query. 25
- Figure 8 Two possible outputs from the same clustering system obtained by changing the required number of clusters with (a) $N = 3$ and (b) $N = 8$. As we can see, the clustering task is a non trivial problem that highly depends on the way parameters are initialized and the level of detail targeted. This parameter selection issue is common to every clustering task, even out of the scope of time series mining. 27
- Figure 9 The three main steps of a classification task. (a) A training set consisting of two pre-labeled classes C_1 and C_2 is entered into the system. The algorithm will first try to learn what the characteristic features distinguishing one class from another are; they are represented here by the class boundaries. (b) An unlabeled dataset is entered into the system that will then try to automatically deduce which class each datapoint belongs to. (c) Each point in the set entered has been assigned to a class. The system can then optionally adapt the classes boundaries. 28
- Figure 10 Example of application of a segmentation system. From (a) usually noisy time series containing a very large number of datapoints, the goal is to find (b) the closest approximation of the input time series with the maximal dimensionality reduction factor without losing any of its essential features. 30
- Figure 11 A typical example of the time series prediction task. (a) The input time series may exhibit a periodical and thus predictable structure. (b) The goal is to forecast a maximum number of upcoming datapoints within a prediction window. (c) The task becomes really hard when it comes to having *recursive prediction*, i.e. the long term prediction of a time series implies reusing the earlier forecast values as inputs in order to go on predicting. 31
- Figure 12 An idealized example of the anomaly detection task. A long time series which exhibits some kind of periodical structure can be modeled thanks to a reduced pattern of "standard" behavior. The goal is thus to find subsequences which does not follow the model and may therefore be considered as anomalies. 32
- Figure 13 The task of motif discovery consists in finding every subsequence that appears recurrently in a longer time series. These subsequences are named motifs. This task exhibits a high combinatorial complexity as several motifs can exist within a single series, motifs can be of various lengths and even overlap. 33
- Figure 14 Complete classification of all the time series representations reviewed in this chapter. 35

Figure 15	Complete classification of the distance measures reviewed in this chapter. 40
Figure 16	<i>Pareto dominance</i> relations for a minimization problem in a bi-criteria space. Any point x of the criteria space divides it into three sub-spaces depending on the dominance relation. $S_{<}$ contains the elements that dominates x ($\forall y \in S_{<}, y \prec x$). Elements of $S_{>}$ are dominated by x ($\forall y \in S_{>}, x \prec y$). Finally, the elements of $S_{?}$ simply cannot be compared to x as they are not dominated nor dominate x ($\forall y \in S_{?}, x \not\prec y \wedge y \not\prec x$). 50
Figure 17	Efficient solutions and dominated solutions for a bi-criteria minimization problem. We can clearly see the <i>Pareto front</i> of non-dominated solutions. The dotted lines define the sub-spaces which are dominated by these solutions. 51
Figure 18	Three efficient solutions x_a, x_b, x_c and their corresponding induced weighted Chebyshev norms N_a, N_b, N_c in a bi-objective problem. Each point of the Pareto front is thus the best solution of a mono-objective problem weighted by its corresponding Chebyshev norm. 52
Figure 19	Illustration of the MOTS matching problem in a bi-objective context. The query Q is at the origin of the space and is represented by a set of time series that have to be matched jointly. Solution \mathcal{A} is the best match for objective \mathcal{O}_1 , as we can see the first time series is closely similar to that of the query. Solution \mathcal{B} is respectively the best match for objective \mathcal{O}_2 . The element \mathcal{C} would be the best solution for the weighted monobjective problem given a set of equal weights. We can see that it is not closely similar to any objective, which motivates the use of multiobjective optimization. 58
Figure 20	Trying to find the solution to the similarity problem exposed in Figure 5 with a multivariate nearest-neighbor approach. The system will order element S_1^1 as being the most similar to S_2^1 , as the distance is slightly different between the two features. Therefore, there is an implicit preference towards the <i>pitch</i> of different sounds in similarity matching. 60
Figure 21	Trying to find the solution of the similarity problem exposed in Figure 5 with a MOTS approach. If we seek to find which elements are more similar to S_2^1 , the system will divide the problem in its two inherent dimensions. Therefore, element S_1^1 and S_2^2 are selected as being the most similar as they are not dominated. Only element S_1^2 is exhibited as being least similar. Therefore, there is no implicit preference towards any dimension in similarity matching. 61
Figure 22	Construction of the quantified bins for time series and computation of the 1st-level distance for a query time series. 62
Figure 23	The approximate lower bounding distances in the criteria space and a set of relationships that can or can not be computed 63
Figure 24	Geometric interpretation of the <i>multiobjective hyperplane search</i> algorithm 66
Figure 25	Query wall time (in seconds) for increasing database size (left) and increasing number of objectives (right) on synthetic datasets. 68

- Figure 26 Space pruning ratio for increasing database size (left) and increasing number of objectives on synthetic datasets. 69
- Figure 27 Query wall time (in seconds) for increasing database size (left) and increasing number of objectives (right) compared between synthetic and real datasets. 69
- Figure 28 Space pruning ratio for increasing database size (left) and increasing number of objectives (right) on synthetic and real datasets. 70
- Figure 29 Shifting from the QBE paradigm (left) to the MOTS framework (right). In the QBE approach (left), a soundfile is fed to the system for which similar sounds have to be found in a database. The system answers with an ordered list of soundfile results. By using time series techniques (center), we can construct a system which match the temporal evolution of any audio feature. However, the combination of multiple audio features input (right) requires a more flexible matching process, hence exhibiting the relevance of the MOTS framework. 72
- Figure 30 Algorithmic framework for two types of interaction. In *MultiObjective Spectral Evolution Query* (MOSEQ), a set of time-evolving properties is drawn. The MOTS algorithm allows to find the set of efficient solutions. In *Query by Vocal Imitation* (QVI), the user can directly use his voice to perform an imitation of the desired properties. A spectral analysis leads to the set of properties. 73
- Figure 31 When a soundfile is input to the system, the analysis module computes a set of descriptors whose mean, deviation and temporal shape are stored separately inside an SQL database. Symbolic information can also be stored in the database, either by automatic extraction or direct user input. 75
- Figure 32 Comparison of different query results for multiobjective optimization and mono-objective selection in a QBE context. (Left) A sound taken from a restaurant scene and belonging to the *crowd* class. (Right) A clip taken from the *female speech* class. 76
- Figure 33 Cross-correlations of the temporal shapes of various audio features analyzed through a hierarchical clustering process. We represent the each cluster with its corresponding information type and representative feature selected. 81
- Figure 34 Experimental interface for the descriptor shape similarity task. Ten sounds are displayed for each of the 21 possible features combination. For each of these sounds, the temporal shapes of two features are displayed. Subjects are asked to rate their perceived similarity between the sound and the temporal shapes. 86
- Figure 35 Interface for the constrained MOSEQ retrieval task. A sound target is presented to the subjects and can be listened repeatedly. The two current features are displayed under drawing boxes. Subjects are asked to draw the temporal evolution of the corresponding features so that it closely match the target. When subjects are satisfied with their drawings, they can perform a query. The MOTS algorithm will provide a set of solutions spread over the optimization space. Subjects can listen to the results, see the temporal evolution of their features and try to find the corresponding target. 88

Figure 36	Interface for the constrained QVI retrieval task. A target sound is presented for each pair of features. The subjects can then perform a vocal imitation of the target. The complete set of sound features are displayed in real-time while the subject is recording. Only the pair of features relevant with the query are marked with a red cross. After recording their imitations, subjects can modify the temporal shapes by direct input. When subjects are satisfied with their input, they can perform a query. The MOTS algorithm will display the corresponding set of tradeoffs solutions. 89
Figure 37	Distribution of skills amongst subjects of the experiment separated between mean distribution of skills (up) and distribution in each of the subjects groups (down) 91
Figure 38	Distributions of mean similarity ratings for each sound feature, independently of the combinations used. Similarity scores for all subjects (up) and group-wise similarity ratings (down). 94
Figure 39	Scatter plots and kernel density estimates for all similarity ratings available and each feature combination. 97
Figure 40	Scatter plots and kernel density estimates for mean similarity ratings produced by each subject. 98
Figure 41	Principal Components Analysis (PCA) of the similarity tournament matrix. 99
Figure 42	Redundancy Analysis (RDA) results performed over the individual score-based directions of listening depending on the groups. 101
Figure 43	Redundancy Analysis (RDA) results performed over the individual score-based directions of listening depending on the subjects. 102
Figure 44	Analysis of the consistency in the directions of listening using a RDA method for the generic similarity task. 105
Figure 45	<i>Similarity heat maps</i> representing the weight of similarity ratings for all subjects and group-wise subjects independently of the features involved. 107
Figure 46	Kernel density estimates of similarity ratings distribution over the optimization space. The estimates takes ratings for all subjects and group-wise subjects independently of the features involved. 108
Figure 47	Mantel spatial correlogram performed over the similarity ratings depending on their distances to each other. 109
Figure 48	Query-dependent analysis of the main <i>effectiveness</i> measures for the MOSEQ constrained retrieval task 111
Figure 49	Query-dependent analysis of the main <i>effectiveness</i> measures for the MOSEQ constrained retrieval task 113
Figure 50	Distributions of mean vocal control strength for each sound feature, independently of the combinations used. Control scores for all subjects (up) and group-wise control (down). 115
Figure 51	Canonical Correlation Analysis (CCA) between the user-wise similarity tournament matrix and the user self-rated skills. 117

- Figure 52 (Left) Hypervolume dominated by a Pareto front given the reference point r_p in a 2-dimensional space. The darker gray subpart defines the box \mathcal{B}_1 which is dominated by point p_1 . The hypervolume \mathcal{H} dominated by the Pareto front is defined as the union of all boxes dominated by each point of the front. (Right) Comparison of two dominated hypervolumes \mathcal{H}_1 and \mathcal{H}_2 . Even though the first class have more elements belong to the final Pareto set, its hypervolume \mathcal{H}_1 is smaller than the hypervolume \mathcal{H}_2 of the second class. 122
- Figure 53 Comparison of distance-based classifiers. The element to be classified is represented by the cross at the origin of the space. The *Nearest-Neighbors* techniques select the class of nearest elements based on the norm of their distance vector. The *Nearest Center* technique first computes the centroid of each class and then selects the nearest one. The *MOTS* paradigm computes the Pareto front and then selects the most represented class. Finally the *HV-MOTS* technique computes the *hypervolume dominated* by each class and then select the class with the largest one. 123
- Figure 54 Comparison of several classification approaches based on the class boundaries that they define in *feature* space. The techniques are separated between their definition of *linear* (left) or *non-linear* (right) class boundaries. 124
- Figure 55 Comparison of the classification boundaries represented in feature space implied by 1-NN selection or HV-MOTS classification algorithms. The first problem (up) represents synthetic data where the classes are *almost* linearly separable. The second problem (down) represents a mixed set of classes data with no linear separation. 128
- Figure 56 Comparison of statistical significance between classification methods based on the Tukey-Kramer HSD over Friedman's ANOVA. 136
- Figure 57 Comparison of statistical significance between classification methods based on the statistical mean difference over a one-way ANOVA. 137
- Figure 58 Comparison of statistical significance between classification methods for an increasing number of classes based on the Tukey-Kramer HSD over Friedman's ANOVA 138
- Figure 59 Comparison of statistical significance between classification methods for an increasing number of samples based on the Tukey-Kramer HSD over Friedman's ANOVA 139
- Figure 60 Comparison of statistical significance between classification methods for an increasing number of features based on the Tukey-Kramer HSD over Friedman's ANOVA 140
- Figure 61 Comparison of statistical significance between classification methods depending on the scientific domain being studied based on the Tukey-Kramer HSD over Friedman's ANOVA 141
- Figure 62 Comparison of the statistical significance for different parameters of warping for the DTW distance as opposed to simple resampling factors of the time series compared with the Euclidean distance. 145

Figure 63	A complete cardiac cycle analyzed through recordings of the heart sounds (PCG), its skin electrical activity (ECG) and pressure in the aortic and atrial valves. 148
Figure 64	Algorithmic workflow for our heart sounds biometry system, summarizing the four milestone of listening which are <i>where</i> , <i>when</i> , <i>what</i> and <i>how</i> to listen. 149
Figure 65	Typical frequency distribution of heart sounds based on the Stockwell transform analysis of 15.814 complete cardiac cycles. The S1 and S2 sounds recorded from 212 different subjects have been processed separately. 150
Figure 66	Illustration of the temporal and frequency resolution of the Stockwell transform. A synthetic signal (left) with very specific properties and its corresponding S-Transform spectrum (right). 151
Figure 67	Computation workflow (up) and specific filter design (down) of the S-Frequency Coefficients (SFC). As we can see, the SFC are computed on a model similar to the MFCC. However, its fundamental differences comes from the use of the S-Transform and its filterbank designed to match the properties of heart sounds. 153
Figure 68	The evaluation of a biometric system in a real-life scenario given its distributions of <i>genuine</i> and <i>impostor</i> scores. 155
Figure 69	Possible tradeoffs between the False Match Rate (FMR) and the False Non Match Rate (FNMR) for different methods 157
Figure 70	Detection Error Trade-off (DET) curves for different methods 158
Figure 71	Receiver-Operator Characteristic (ROC) curve 160
Figure 72	Results of the system identification based on the Rank-k identification rates 161
Figure 73	Results of the menagerie analysis performed over every heart beats for the best feature combination. 162
Figure 74	Comparison of the classification accuracy of using <i>only</i> temporal features, <i>only</i> static (mean and deviation) features or <i>mixed</i> sets of information with either multiobjective or mono-objective selection. 171
Figure 75	Temporal profiles drawn by each of the 19 participants for the set of 6 morphological classes. 175
Figure 76	Accuracy of different classification methods for an increasing number of descriptors over the sound morphology. 176
Figure 77	Results of a hierarchical clustering performed on the sounds of the study using the audio features selected thanks to the HV-MOTS classification paradigm. 177
Figure 78	Original approaches to tackle the problem of computer-aided orchestration. A target sound file is fed to the system which will try to reconstruct it by using a set of instruments. The system rely on a set of feature files which are combined through prediction functions. 184

Figure 79	A new approach for computer-aided orchestration. Unlike previous systems (cf. Figure 78), the comparison between targets and sound mixture is based on their temporal evolutions. The knowledge is encapsulated in a SQL database in order to provide an almost infinite source of knowledge. Finally, the need of a well-formed audio example is bypassed by the direct input of temporal shapes.	186
Figure 80		188
Figure 81		190
Figure 82	We combine paradigms by using the MOTS algorithm in order to find efficient solo instruments that can be used as “seeds” to the initial population. That way we reduce the blindness of pure randomness	191
Figure 83	The current prototype for Abstract Temporal Orchestration with Modular Structure (ATO-MS) features an extensible architecture of modules to tackle the problem of Computer-Aided Orchestration.	193
Figure 84	Make an instance of the db-sound class by a SQL query. The new db-sound instance contains four sounds. A special editor displays the current one and different descriptors of it.	197
Figure 85	Concrete example of the usage of spectral maquettes. Several functional, macro and micro-temporal relations are defined between boxes.	199
Figure 86	A complete system of constraint inference that allow to bridge the gap between the symbolic realm of musical writing and the signal world of timbre. A simultaneous analysis of the symbolic score and corresponding audio features evolution could provide a graph of constrained relationships, explaining the link on both types of viewpoints.	209
Figure 87	Comparing the resolution power of the FFT to the S-Transform for a single cardiac cycle.	240
Figure 88	Influence of the despiking filter exhibited through the ROC curves of <i>with</i> and <i>without</i> use of the filter.	244
Figure 89	Influence of the wavelet denoising exhibited through the ROC curves of <i>Coiflets</i> and <i>Daubechies</i> wavelets.	245
Figure 90	Influence of the segmentation parameters exhibited through the ROC curves of spectral differences (F+ or F-), window size and number of partials.	247
Figure 91	Influence of the beat selection exhibited through the ROC curves of <i>energy deviation</i> and <i>shape deviation</i> .	248
Figure 92	Influence of the size of the warping window for comparing the time series with the DTW distance measure exhibited through the ROC curves of 2, 5, 10 and 20% of authorized warping reach.	249
Figure 93	Influence of the size of the resampling factors for comparing the time series exhibited through the ROC curves of 128, 64, 32, 16 and 8 resampling points.	251
Figure 94	Influence of the type of <i>decision rule</i> exhibited through the ROC curves of the <i>mean</i> , <i>min</i> and <i>max</i> rules.	252
Figure 95	Influence of the size of <i>testing set</i> exhibited through the ROC curves of different set cardinalities.	253

Figure 96	Influence of the size of <i>training set</i> exhibited through the ROC curves of different set cardinalities.	254
-----------	---	-----

LIST OF TABLES

Table 1	Comparison of the distance measures surveyed in this chapter with the four properties of robustness. Each distance measure is thus distinguished as <i>scale</i> (amplitude), <i>warp</i> (time), <i>noise</i> or <i>outliers</i> robust. The next column shows whether the proposed distance is a metric. The cost is given as a simplified factor of computational complexity. The last column gives the minimum number of parameters setting required by the distance measure.	44
Table 2	List of available descriptors whose mean, deviation, temporal shape and first and second derivatives are stored separately. More detailed information can be found in [295]	74
Table 3	Selected features and information class based on the analysis of cross-correlations between sound features of the dataset used in the reminder of this study. The last column contains the corresponding set of features that are strongly correlated to the selected one.	83
Table 4	Tournament-based analysis of listening directions for the <i>shape similarity</i> task.	95
Table 5	Linear coefficients of the PCA for the three main components displayed in Figure 41	99
Table 6	Tournament-based analysis of listening directions for the generic similarity task.	103
Table 7	Tournament-based analysis of listening directions for the constrained retrieval task.	106
Table 8	Overall similarity ratings for the pertinency task depending on the feature and the users groups.	109
Table 9	Results of the effectiveness of the MOSEQ paradigm evaluated through the <i>mean</i> and <i>cumulative</i> statistics for the <i>task completion time</i> , <i>number of queries required</i> and <i>number of features modifications</i> required as well as the <i>mean number of audio files played</i> and the corresponding <i>time required</i> to play these files. The last part of this table displays the <i>feature-dependent mean task completion time</i> .	111
Table 10	Results of the effectiveness of the QVI paradigm, evaluated through the <i>mean</i> and <i>cumulative</i> statistics for the <i>task completion time</i> , <i>number of queries required</i> and <i>number of features modifications</i> required as well as the <i>mean number of audio files played</i> and the corresponding <i>time required</i> to play these files. The last part of this table displays the <i>descriptor-dependent mean task completion time</i> .	112
Table 11	Results of user satisfaction survey for the MOSEQ constrained retrieval task	113

Table 12	Results of user satisfaction survey for the QVI constrained retrieval task.	114
Table 13	Tournament-based analysis of listening directions for the constrained retrieval task.	116
Table 14	Comparison of overall classification accuracies for different methods	135
Table 15	Comparison of classification accuracies with state-of-art results on the same datasets. We provide for each dataset the original algorithm used to obtained the reported classification accuracy.	143
Table 16	Comparison of classification accuracies with state-of-art results on the same datasets. We provide for each dataset the original algorithm used to obtained the reported classification accuracy.	144
Table 17	Details of the PCG recordings datasets. We provide the	154
Table 18	Result of <i>Order-0 analysis</i> for different methods	156
Table 19	Result of <i>Order-0 analysis</i> for different levels using the HV-MOTS method	156
Table 20	Result of <i>Order-1 analysis</i> for different methods	157
Table 21	Result of <i>Order-1 analysis</i> for different levels using the HV-MOTS method	158
Table 22	Result of <i>Order-2 analysis</i> for different methods	159
Table 23	Result of <i>Order-2 analysis</i> for different levels using the HV-MOTS method	159
Table 24	Result of <i>Order-3 analysis</i> for different methods	159
Table 25	Result of <i>Order-3 analysis</i> for different levels using the HV-MOTS method	160
Table 26	Results of the analysis of the <i>template ageing</i> phenomenon.	163
Table 27	Lists of both user and environmental factors that could potentially affect the performances of the heart sound identification system.	164
Table 28	Description of the MuscleFish dataset used in classification tasks. 409 sounds are divided into 16 classes.	168
Table 29	Description of the Freesound dataset collected specifically for our study. 2193 sounds are divided into 54 classes.	169
Table 30	Classification results on the MuscleFish dataset for a growing number of objectives. For a given number of objectives, the left column indicates the mean classification accuracy and the right column indicates the best classification accuracy.	169
Table 31	Classification results on the Freesound dataset for a growing number of objectives.	170
Table 32	Significance tests between various methods for a growing number of objectives across both datasets. For a given number of objectives, the left column indicates the mean column rank (from Tukey-Kramer HSD over Friedman's ANOVA) and the right column gives the statistical mean difference in accuracy with the top performing method (from a one-way ANOVA).	170
Table 33	Confusion matrix for the best classification accuracy (95.4%) obtained by HV-MOTS on the MuscleFish dataset. The descriptor combination used is composed of MFCC, MFCCDeltaStdDev, PerceptualSlope, ChromaDeltaStdDev, RelativeSpecificLoudnessDeltaStdDev and PerceptualDecrease.	172

Table 34	Effects of a set of distortions on classification accuracy for different methods on the MuscleFish dataset. 174
Table 35	The best combination (6 features) obtained thanks to the HV-MOTS classification paradigm provides a classification accuracy of 87.3%. The right column shows the individual classification accuracy for each feature. 177
Table 36	Comparison of algorithms on the monophonic orchestration problems 192
Table 37	Comparison of algorithms on the polyphonic orchestration problems 192
Table 38	Comparison of the orchestral databases used as a knowledge source for <i>Orchidée</i> (left) and our <i>ATO-MS</i> system (right). 195

ACRONYMS

DRY	Don't Repeat Yourself
API	Application Programming Interface
UML	Unified Modeling Language

Part I

INTRODUCTION

The onset of this study takes its roots in musical orchestration. Orchestration is the subtle art of writing musical pieces for the orchestra, blending the sounds of diverse instruments together by taking into account the acoustic specificities unique to each. The origin of this art lies in the willingness of composers to empower the orchestra to become an expressive unity that can accurately transcribe emotions. This goal can be reached through the knowledge and educated use of the different spectral qualities of each instrument. Composers can then build and adjust particular emotional effects over time. Piston [303] wrote in his treatise that orchestration *"is aimed at discovering how the orchestra is used to translate a musical thought. It is a means to study how the instruments blend together to create equilibrium of sounds"*. If we consider the range of expressivity offered by a single instrument, we can get a glimpse on the extent of sonic possibilities offered by an orchestra. Furthermore, given the variety of notes, playing modes and dynamics that can be obtained by an instrument, we can clearly see the combinatorial complexity that is embedded in the art of musical orchestration. Beyond its traditional sense, orchestration rely heavily on the concept of sound mixtures, which ubiquity range nowadays from orchestral to electronic music. Orchestration can be thought as the realm of musical writing in which the timbre acts as the main parameter (we shall try to provide a definition of the timbre in Section 2.3). Amongst all the components of music writing, orchestration has long remained in his teaching as in its practice, an empirical activity. Only quite recently has risen the idea of computer-aided orchestration, which is faced with the vastness of its field of study. Indeed, the topic of orchestration encompasses notions from auditory perception, music analysis, music theory, composition, signal processing and computer science. The difficulty of finding a rigorous formalism along with its youth in the musical discourse still makes orchestration one of the ultimate musical dimensions that have not been studied enough to fit its complexity. We must, therefore, start by investigating the reasons behind this paucity of researches by looking into the past of orchestration in order to envision its future.

1.1 FROM THE PAST EMPIRICISM ...

In his earliest essay, Berlioz [40] laid down the foundations of what would later become the teaching of musical orchestration. Already at his time, he noted that the art of orchestration is *"taught as little as the ability to find beautiful songs, beautiful successions of chords and original and powerful rhythmic forms"*. Even if we can clearly define a song, a succession of chords or a rhythm, it appears unfeasible to interpret objectively the beautiful nor the original. A few number of treatises followed his work, ongoing to the almost encyclopedic work of Koechlin [222]. From these seminal studies to contemporary works such as Piston [303] and Casella [74], the problem of orchestration is still always exposed and taught to composers on an empirical basis. As there are no fixed, generic rules, such as those that govern harmony in Western music, the exploration of combinatorial possibilities offered by instrumental properties is always reduced to a series of "orchestral recipes". Thus, even in the most recent works, we can

learn that the flute is an instrument with typical pastoral shades or that the bassoon produces tones that resemble a libidinous old man. As strict, universal orchestration rules continuously seem to have eluded scholars over several generations, the authors of treatises delve in producing a collection of examples drawn from their own experience and explore the resulting orchestral color in their own terms. The various treaties thus boil down to a series of recipes identifying some of the orchestral archetypes. It is clear that these empirical approaches of orchestration preclude the systemic use of such treatises on a scientific basis.

In his practice of orchestration, the composer must face the choice of using his own (necessarily limited) personal experience or to turn to orchestration treatises that fail to examine the combinatorial possibilities of orchestral timbre on a rigorous basis. As full of examples as might be all the treatises, the question of their accuracy and scope, however, deserves to be asked. First, as exposed before, treatises lack the exploratory component that could provide a structural support on instrumental mixtures for the composers to adjust the final outcomes. Furthermore, recent developments in instrumental music have focused on the introduction of exotic playing modes that allow instruments to produce sounds previously unheard. Therefore, as most of the orchestration studies date back to several decades, there seems to be a potential obsolescence of their repertoire of the study. The language used by various authors may also already reflect the aesthetic vision of their time. From all these observations, it does not appear that an orchestration treatise on rational grounds may exist to this day; however, this need seems crucial.

1.2 ... THROUGH MUSICAL WRITING ...

Even if the art of orchestration focus on the use and evolution of sound properties, it remains an act of musical writing. In fact, it is extremely difficult to separate these two aspects as a musical material is often thought and written for a pre-defined set of instruments. However, orchestration cannot be enclosed in a single frozen time of writing, or limited to the field of instrumental technicalities. Apart from isolated chord situations where the pitches of all instruments are stable, every other musical context (figures, textures, gestures) seem to fall together within a tied writing and orchestration. In its daily practice of orchestration, the composer does not only superimpose stationary sounds but instead focus on carving, vertically and horizontally the sound material. He can simulate the attack of an instrument by another, with a third one acting as resonance. The position of the orchestration in relation to the entire musical knowledge have significantly evolved through different eras and authors. Koechlin, for example, seems to consider it as a technique totally subservient to the other dimensions of writing: "*We must carefully plan every orchestral element such as accents, cadences, progressions and dynamic overlaps in relation to a piece as a whole. Each change should be based on its musical context*" [222]. From these observations, we can distinguish two modes of writing for the orchestra. First, a "holistic" orchestral writing where melodies and instruments are thought of as indivisible units, which we call *inductive orchestration*. On the other hand, an "abstract" writing in which the score and chords are produced independently and then orchestrated, which we call *projective orchestration*. Koechlin advised against this practice, stating "*the too much prevalent belief that the orchestration is the basic allotment of timbres in various instrumental lines is clearly inadequate*" [222].

- *Inductive orchestration* is created by having precise orchestral colors and effects in mind and trying to produce combinations of instruments that could achieve these ideas.
- *Projective orchestration* is produced by first writing an “abstract” score and then determining the allocation of different melodic lines to instruments. Although Koechlin decried this practice, famous examples include Maurice Ravel who orchestrated his piano plays and, therefore, had not thought about the instrumental colors when writing his original pieces.

From these two potential forms of orchestration, composers seem to be more keen on inductive orchestration, thus writing musical pieces with a precise idea of their orchestral colors. In both cases, the art of orchestration is directly related and even almost indivisible from musical writing, as they are not isolated processes from the compositional practice. Once again, the problem of music writing is an intricate area for scientific research that encompass techniques that fall under subjective assessments, and whose mechanisms appear difficult to formalize.

1.3 ... TOWARDS A MODERN TREATISE ?

Given this empirical tradition, it seems almost hazardous to tackle the orchestration on a rigorous scientific basis. However, Bregman [58] envisioned, “*it should be possible to write an orchestration treatise based on fairly abstract principles to be applied to all styles of music. This would not comprise a collection of precepts accumulating imperatives and prohibitions, but a guide to how to get a particular sound, leaving the composer free to alter the outcomes*”. When Bregman refers to *all* styles of music, we can see that the breadth of musical orchestration goes beyond instrumental writing and can be generalized to the notion of sound mixtures, mainly when considering the latest trends in contemporary music.

The advent of electroacoustic music has produced a fundamental shift in instrumental music, which seems now to evade from the traditional categories of musical writing (*pitch, duration, mode*) and heads itself towards the composition of inharmonic and noisy sounds which exhibit large timbre variations over time. The most evident marks of this evolution are the *spectral music* (Grisey, Murail), *concrete instrumental music* (Lachenmann, Sciarrino) and the *school of complexity* (Ferneyhough, Dillon) which seek to “*saturate writing*”. It is interesting to note that this problem seems specific to Western music, locked in his writing since the Gregorian era. Indeed, the music from other cultures, mostly passed down through oral traditions has long been attached to address these categories of complex and noisy musical elements. Amongst nowadays composers, many are focusing more resolutely towards the expressive qualities of sound and the potential of complex sounds and noises offered by the infinite capabilities of electronic source materials. Musical orchestration in the context of these sonic possibilities becomes an even more complex operation. This question brings us back to the fuzzy boundaries that exist between sound and noise. With the ceaseless flourishing of electronic music, it seems that this limit has been pushed further into a corner. The advent of technology has made us consider noisy sounds as potential units of musical creativity, evidenced by the approach of Pierre Schaeffer [340]. He chose to compose with sounds and noises recorded from any potential source in the every day environment, thus giving rise to the *musique concrète* movement. He exposed the need to “*replace the limited variety of instruments that constitute an orchestra with the infinite variety of timbres provided by noises obtained through special mechanisms*”. Schaeffer suggested

elevating noises to music by transforming and repeating such sounds in order that the listeners' awareness might be directed by temporal expectations. This approach thus seeks to create a causal dependence between elements in order to delineate a clear boundary between music and noise. These methodologies show us that the process of organizing sound properties can be highly conditioned by the elementary units that are sought to be combined.

1.4 ELEMENTARY UNITS OF ORCHESTRAL STUDY

As discussed previously, the art of musical orchestration embeds all the techniques that attempt to combine musical elements and "*consider how they can contribute to an emotional effect in their musical association*" [40]. In the classical orchestra, we can perform a categorization of these elements on several scales, starting from the instruments (violin, cello, clarinet, flutes), the families of instruments (strings, woodwind, brass) and up to the mode of production. However, even inside a common production apparatus, the sound characteristics of each instrument provide an almost infinite variety. A tremendous scope of expressivity can be found in each instrumental repertoire of playing modes, range of notes, intervals and chords, range of dynamics and finally the overall structure and rhythm that creates a musical context. Even with all these parameters fixed, the nature of instrumental sounds still remains somehow stochastic. Thus, we must keep in mind that even at the same pitch, dynamics and duration, the resulting signal can vary from one instance to another, which further deepens the combinatorial complexity of the orchestral study. This variability of elements appears through different dimensions and can be observed on multiple temporal scales. Many of those are still to this day oblivious to scientific study. Nevertheless, in order to combine these elements appropriately, we need to consider their range, playing modes, playability, dynamics, articulations, speed, role and purpose in an ensemble and even their individual ability to translate musical thoughts. As underlined by Piston [303], "*we can not overstate the importance of instrumental knowledge. An appropriate instrumental writing is undoubtedly the determining factor in the success of an orchestration*". Thus, it seems essential to take steps towards the support of the modern instrumentarium in which the computer seems to have an evergrowing importance. If the musical research does not capture these individual sound particularities as an aim in itself, the orchestration will remain an empiric and haphazard science. As we try to devise a formalism to pave the way for modern orchestration, we thus need to go through with a deep understanding of the sound properties of its elementary units.

2

STUDYING ELEMENTS

In the study of musical orchestration, a fundamental question lies in the selection of elements to combine. Indeed, how to integrate the properties of different sounds, if we are not yet able to master these constituents and how they are perceived.

2.1 MUSICAL ATOMS

The study of orchestral pieces can be performed on several distinct scales as exemplified in Figure 1. The infinite variety of instrumental combinations can generate widely varying sound results, therefore, inducing an aspect of *combinatorial complexity*. At the same time, the temporal evolution of musical structures strongly shape the perception of musical pieces. These *macro-temporal structures* can create *tension, expectation, release* and a wide variety of emotions depending on their construction. However, if we focus our attention to even a *single* element inside these complex structures (such as the red-circled note in Figure 1), another wealth of complexity reveals itself. These elements, that we will call *musical atoms*, also exhibit unique temporal evolutions of their spectral properties. These *micro-temporal properties* vary even for the same instrument at the same pitch and loudness, as we shall discuss in Section 2.3. Furthermore, the temporal evolution of multiple decorrelated sound properties can be perceived simultaneously leading to a *multidimensional perception*, which we discuss in details in Section 3.3.

We call these elements *musical atoms* in the broadest sense that they constitute a coherent spectral unit. The first element that should come to mind is a single note from an instrument, but a continuous glissando also form a logical unit. Therefore, we also embed in this definition every acoustic entity that represent a *unit of musical creativity*. Indeed, as exemplified by the approach of Pierre Schaeffer [340], even everyday noises can be elevated to the rank of musical units depending on the temporal organization of such elements.

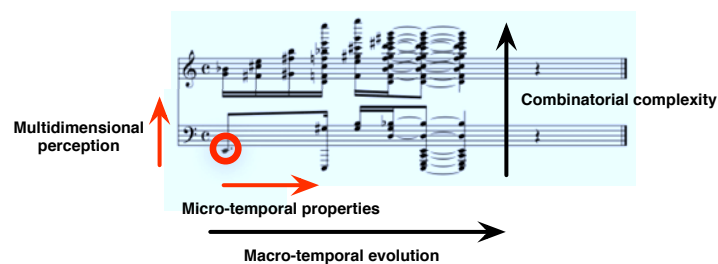


Figure 1: Different levels of complexity appear in the study of music, where the processing of musical pieces is performed on several scales simultaneously. The various potential combinations of instruments induce a *combinatorial complexity*. The *macro-temporal evolution* of these musical structures strongly condition our perception. However, if we focus on a single musical atom, it also embeds *micro-temporal* evolutions of its spectral properties that are perceived in a *multidimensional* fashion.

2.2 SIGNAL AND SYMBOLISM

Composition may be seen as the act of projecting a hypothetical signal that exists only in a composer's mind to an efficient symbolic representation. This strategy has been used for centuries as it allows sharing this mental image to performers that will try to reproduce it in the most accurate way. Western music was largely based on this *harmonic paradigm*, i.e. it built its musical concepts from sounds considered harmonic and stationary. Western musical notation is also without coincidence, consistent with this paradigm (writing a pitch as a point is reducing a sound to its fundamental, and considering that, inside this sound, changes in the harmonics follow a parallel and constant evolution). However, this symbolic tradition has now to coexist with a consideration of composers increasingly marked with the spectral qualities of musical elements. The musical evolutions have led us to reconsider the inharmonic and noisy sounds as part of the musical perception. A substantial theoretical debate in the study of sound mixture composition comes from the use of multiscale representations for linking heterogeneous data (cf. Figure 1). From this, unfold a contemporary issue in computer music research : the *signal / symbolic interaction*. These two research streams have long remained impervious to each other, partly because of the apparent heterogeneity of their objects of study. On one hand, the analysis and synthesis of digital signals contributed to the comprehension and production of sounds previously unheard. On the other hand, the symbolic approach focused on the analysis of musical notation structures, but the composition of sound mixtures is located precisely at the intersection of these two lines of research. If it claims to produce timbres, it is through a symbolic writing process. Therefore, it is an ideal meeting point between the symbolic and spectral domains. Sound as a material coexists and interacts with formal structures. Nowadays, computers offer the possibility to manipulate a sound object in a compositional process, while simultaneously studying its acoustic properties through analytical tools. We can thus combine the symbolic discipline of writing with the spectral domain of timbres possibilities. It should, therefore, become possible for the composer to relate the exploration of sound to the organization of symbolic data.

2.3 MULTIPLICITIES OF TIMBRE

As stated earlier, the *timbre* is the dominant parameter of this study, as musical orchestration is focused towards attaining new timbre through the combination of individual instrumental timbres. However, we carefully avoided defining and overusing this term, and preferably used the instrumental *sound properties* or *qualities*. Indeed, when the musical timbre is put forward in a scientific study, the first problem to arise is that of its precise meaning. It seems that a consensual definition of musical timbre has eluded scholars for several generations. Attempting to formulate a unique and accurate definition inevitably leads to a number of difficulties, as the systematic use of this word has finally made it airtight and elusive. Despite a century of consecutive studies, there seems to be no consensus on a comprehensive definition from which researchers could build models or theoretical methods. The timbre rather encompass a multitude of definitions that are formulated differently by acousticians, musicians or computer scientists, as we will detail in the remainder of this section.

The American Standards Association (ASA) defines timbre as "*that attribute of sensation in terms of which a listener can judge that two sounds having the same loudness and pitch are dissimilar*" [17]. From this definition, the timbre would seem like a property

that could intervene only in discrimination tasks but can not be quantified positively. In simpler terms, we could say that the timbre of an instrument is its collection of sound *characteristics* which allows to recognize it and makes it differentiable from other instruments even at the same pitch and loudness.

2.3.1 *Timbre and acousticians*

Following the work of Joseph Fourier on the decomposition of periodic functions in 1822, the first significant research on the instrumental timbre can be traced back to the seminal work of the physicist Hermann Helmholtz [175] in 1863. At this time limited by his means of investigation, he is forced to restrict his analysis to sustained sounds, which he called *musical tone*. By using a series of resonators, he evidenced the harmonics of instrumental sounds and then suggested that the timbre is defined by the mean intensity of the harmonic components. However, he already noted "*with a little thought we also see now that some of these acoustic features depends on how the sounds start and finish*". Following the work of Helmholtz, the study of instrumental sounds has long remained entirely ignorant of their temporal aspects. This approach is now strongly avoided (as we will discuss further in Section 2.4) as it is widely recognized that the perception of timbre can not be separated from its temporal aspects. Following the work of Risset [329], it is now well-accepted that the characteristics of timbre are heavily influenced by the temporal variations of *each* harmonic component. His study also revealed the dynamic nature of the energy envelope and the importance of the attack segment in recognition of sounds. Therefore, we can define the *causal timbre* that is involved in recognition of the sound source. On the other side lies the sound properties that provide a qualitative perception of sound. The timbre is thus both the identity of the sound source and sound qualities it possesses.

2.3.2 *Timbre and musicians*

It seems that the appearance of the timbre in the musical discourse comes from the work of Berlioz [40]. He introduced his treatise by saying : "*The purpose of this book is first, an indication of the extent of the instrumental mechanisms. Then we shall turn to the so far neglected nature of the timbre, the nature and expressive abilities of each instrument and the best known methods for grouping them properly*". Thus, when defining timbre, the composers take into account the character of each instrument along with its expressive capabilities. In the history of Western music, the instruments that were originally designed as substitutes for the human voice, have progressively become "*voices*" on their own, whose identity favored polyphonic sound listening. The instruments then started to be employed as entities providing variations of colors in musical discourse. In this musical context, it appears that the two complementary aspects of the timbre are still an integral part of musical creation: the *identity* of a recognizable part and the *evolutions* of its colors. Furthermore, when considered in terms of music writing, timbre plays a dual role at the two poles of the instrumental universe which are the *articulation* and *fusion* [58]. The *articulation* of music is created by the temporal evolution of its structure, based on the possible identification of individual timbres. On the other hand, the *fusion* is the effect that can make it impossible to determine the individual components of a sound mixture. It is thus intended to lead listeners towards new elements of the musical discourse. Until the 19th century, the timbre seems to have played an almost decorative role in musical composition. Romantic composers such

as Hector Berlioz began to consider the timbre as part of the expressive power of music. Arnold Schoenberg and the Vienna School tried to expand its influence in the music discourse. Contemporary composers now use the timbre as a central element of musical aesthetics. Unfortunately, despite their interest in the timbre, composers might not take advantage of its ubiquity, as the timbre is created through complicated physical phenomena which are, therefore, difficult to describe and use.

2.3.3 *Timbre and computer scientists*

Through all the potential definitions of timbre, the notions of sound *qualities* or sound *characteristics* are always prevalent. Therefore, over the past decade, a tremendous amount of work has been devoted to finding relevant high-level features to compute over sound signals. These researches have been made possible by advances in computation speed of the Fourier transform through the Fast Fourier Transform (FFT) algorithm. Over the years, an increasing number of audio features have been developed in order to provide a characterization of the timbre by decomposing it into a set of complementary features. These can be coarsely divided into six main categories; ie. *energy*, *frequency*, *harmonic* (computed only on the harmonic peaks), *spectral* (computed on the whole distribution of the spectrum), *noise* (computed after removing the harmonic peaks) and *perceptual* (computed after filtering the signal with a model of the human ear). Each of these features represents the temporal evolution of a particular characteristic of the related sound. It has now become straightforward to compute this whole set [295] over a sound signal in order to obtain different aspects of the sound. Psychoacoustic studies have further shown the correlation between these sound features and perceptual dimensions. One of the best-known example is the *spectral centroid* that has been shown to be related to the *brightness* of sounds [189]. It has also been shown that the *attack time* heavily influence the perception of percussive aspects of a sound [150]. Other perceptual dimensions have been studied like *sharpness* that relates to the position of *harmonic energy* in the spectrum [392] or *tension* that seems to be primarily connected to *roughness* [309]. Hence, in the computer scientist approach, the timbre is sought to be fully characterized by a set of audio features as comprehensive as possible.

2.4 TIME SCALES CONTINUUM

From a purely aesthetic perspective, it appears obvious that a collection of beautiful things placed upon another with no structure is clearly inadequate. Thus, music can not be represented by a collection of "beautiful" instants frozen in time but is rather developed through a carefully planned temporal structure and progression. Once again, the approach of Schaeffer has shown us that a collection of unpleasant elements can be transformed into music if they are given the proper temporal logic. Every composer knows that the sound of a note depends on its context and that it is the sound pattern, not the isolated note, which determines the auditory perception.

The previously stated definition of timbre by the ASA in 1960 [17] is accompanied by a note (p.45) which states, "*timbre depends primarily on the spectrum of the stimulus, but it also depends on the waveform, sound pressure, spatial position and temporal characteristics of this stimulus*". It is now well-accepted that the timbre does not only depend on the average spectral distribution of sound spectra, but is also strongly tied to their temporal characteristics. Risset [329] found that the proportion of harmonics in the spectrum of brass instruments increases with their loudness. This shows that the timbre can not

be described by the static values of many parameters, but is instead characterized by the temporal evolution of some of these parameters. This is also confirmed by several works on the perception of timbre [57, 275, 342] that demonstrate the importance and primacy of temporal descriptions. As argued previously, the spectral characteristics of timbre are dynamic and continuously evolving elements.

From all these observations, it appears that time should be the primary topic of research in any musical approach, even at the smallest temporal scales. If orchestration requires tools for the vertical and horizontal arrangement of sounds, we have to formalize processes that evolve over time. The revolution of the twentieth century by the emergence of the timbre and discovery of other cultures urge to overpass in one way or another, the paradox of writing for harmonic instruments. In order to accompany this aesthetic paradigm shift, it seems essential to better understand, describe, analyze and compare complex sounds by firstly taking into account the temporal aspects of timbre.

PERCEIVING ELEMENTS

In order to study the potential approaches to combine musical atoms we have seen that we must first understand what we can and can not perceive, but most of all *how* we perceive them. All along this discussion, we will take an analogy with the visual perception in order to clarify our ideas to the readers.

3.1 THE DOORS OF PERCEPTION

The physical perception of living beings has been shaped through millions of years of genetic evolution. This physiological shaping is based on a response to survival needs and for the beings to be able to understand, identify and interact with their environments. The continuous and consistent modifications of the sensory systems can be seen as an indirect action, in which the living species are confronted to various environments to which they should adapt to survive. Hence, what we perceive is based on the evolution of our needs in our environment. Therefore, all forms of perception are meant to provide an *organization* and possible *interpretation* of the surroundings. If we look at the limits of the visual perception, presented in Figure 2, it seems that our visual system is limited to an extremely narrow segment of the electromagnetic spectrum. However, in this small fraction of sensory information that we call *visible light*, we can still distinguish millions of colors and perceive extremely complex and detailed spatial structures.

We now turn our attention to the limits of our auditory perception, presented in Figure 3. As for visual perception, we can only perceive a very small fraction of the acoustic spectrum. However, in this narrow portion of the air vibrations, we are able to differentiate hundreds of pitches and perceive complex temporal evolutions and structures. Finally, even if we consider a fixed part of this acoustic spectrum (supposedly here instruments at the same fixed pitch), we are still able to differentiate between all these sources based on their various properties. Hence, to take a computer scientists analogy, we could say that even if our hardware is somehow limited, our software seems pretty efficient for making the most of this restricted information. Therefore, our sensory systems are meant to provide an *organization* of our surroundings, but it is the *interpretation* of this sensory information that enable our perceptual systems to achieve *identification* of their constituents. Therefore, perceiving elements is being able to identify elements. In turn, in order to be able to perform these distinctions between elements, we need to rely on a complex assessment of their (dis)similarities.

3.2 COMPLEX AND MULTIFACETED SIMILARITIES

3.2.1 Assessing visual similarity

We start by taking a visual example in order to make explicit the difficulties in defining the notion of similarity, that we will subsequently apply to sound perception. The analysis of visual similarity is presented in Figure 4. If we consider a set of three basic visual units (*square*, *circle* and *triangle*) and use only the three elementary colors (*red*,

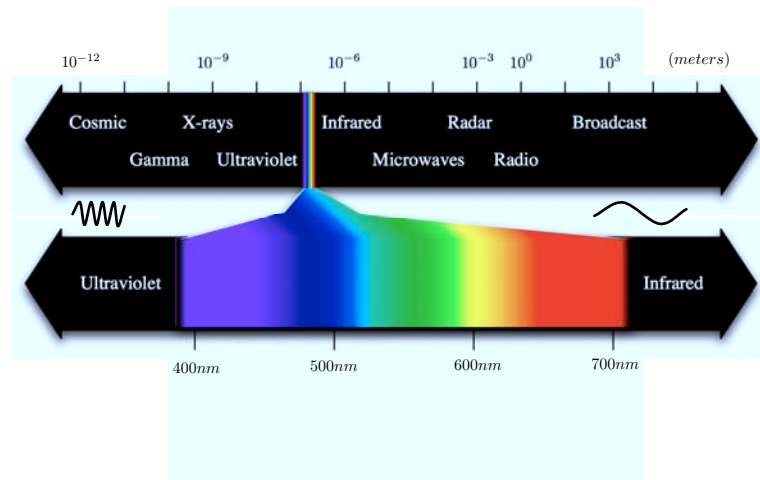


Figure 2: The doors of our visual perception. We are able to perceive only a very small fraction of the electromagnetic spectrum. However, even in this narrow perceivable part we can still differentiate between millions of colors and perceive extremely complex and detailed spatial structures.

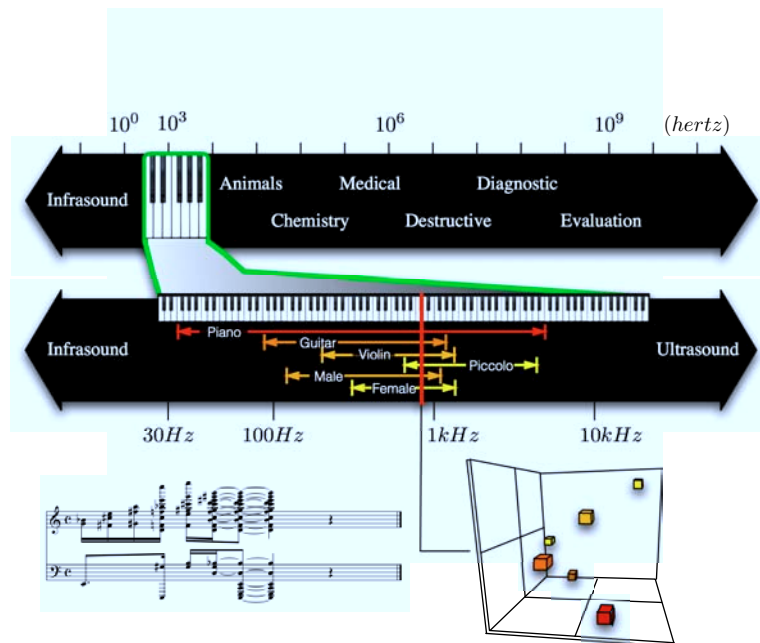


Figure 3: The doors of our auditory perception. As for visual perception, we can only perceive a very small fraction of the acoustic spectrum. However, in this narrow portion of the air vibrations, we are able to differentiate hundreds of pitches and perceive complex temporal evolutions and structures. Finally, even if we take a fixed portion of this acoustic spectrum (supposedly instruments at the same fixed pitch), we are still able to differentiate between all these sources based on their various properties.

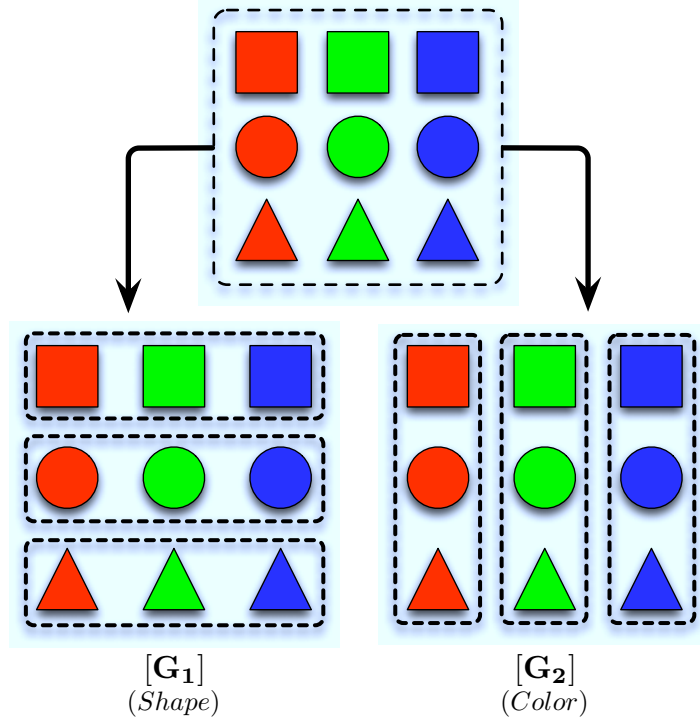


Figure 4: The assessment of *visual* similarity in a set of elementary units. Given this set of elements, we are faced between evaluating their similarity based on the *shapes* of elements such as grouping G_1 or rather based on their *colors* such as grouping G_2 .

green and blue), we can obtain a set of nine elementary units. Given this set of elements, we are faced between two forms of similarity assessments. We can either process the similarity based on the *shapes* of elements such as grouping G_1 or rather based on their *colors* such as grouping G_2 . Being aware of these two features when trying to process the similarity makes it almost impossible to provide a *single* measure of similarity between these elements. For example, if we try to rank the similarity between the red circle and all other elements, it seems unfeasible to determine if the most similar item is one of the circles or another shape that is filled with the same color. This single similarity score (ie. finding *the* most similar element) can not be provided without putting some *preferences* over the two dimensions of variability (*shape* or *color*).

However, this case already allows to get a first sight on what could be the notions of *subjective* and *context-dependent* similarity ratings. Indeed, if this document has been printed in black-and white (and same goes for color-blind people reading it), then the G_2 grouping becomes less relevant, and the similarity can be processed solely through the shape of objects. Therefore, the similarity between elements might be assessed following different *directions* that can vary depending on the subject and the context.

3.2.2 Assessing sound similarity

Now let us turn our attention to the assessment of sound similarity. In order to do so, we consider the most elementary sound that could ever be synthesized or heard,

namely a *sinusoidal signal*. This primary unit that forms the basis of all harmonic sounds can be synthesized by computing the temporal signal

$$x(t) = A(t) \cdot \sin(2\pi\mathcal{F}_0(t) \cdot t + \phi_0) \quad (3.1)$$

with $A(t)$ the *amplitude* function over time (the “loudness”) and $\mathcal{F}_0(t)$ the *frequency* function of the sinusoid (its “pitch”). In our everyday life, we can easily differentiate the pitch and loudness of sounds, so it is fairly painless imagining a sound with an increasing pitch and decreasing loudness simultaneously (for instance imagine a squeaking door or an ascending whistle moving away from you). So we are here studying the simplest element that we could ever find in audio processing, through its most basic properties. A graphical interpretation of this assessment is presented in Figure 5. We can easily define two different temporal functions for the *loudness* and *pitch* properties. We suppose that the temporal evolution of each feature can either be *descending* or *ascending*. Therefore, the temporal evolution of the amplitude $A(t)$ can either be A^1 or A^2 and the temporal evolution of the frequency $\mathcal{F}_0(t)$ can either be \mathcal{F}_0^1 or \mathcal{F}_0^2 . Given these two possibilities for each feature, we can synthesize the set $\mathcal{S}_{\mathcal{A}}^{\mathcal{F}_0}$ of four different sounds s_1^1, s_2^1, s_1^2 and s_2^2 . Now if we try to determine the similarity between elements inside this set, we face the same problematic as for visual similarity. We can either evaluate the similarity between elements based on the evolution of their *loudness* properties such as grouping G_1 or rather based on the temporal evolution of their *pitch* such as grouping G_2 . As discussed in the previous section, it seems impossible to provide a *single* measure of similarity without imposing some *preferences* over the two dimensions of variability.

3.3 MULTIDIMENSIONALITY OF TIMBRE PERCEPTION

As we discussed earlier, several psychoacoustic studies have shown the correlation between sound features and perceptual dimensions. The previous example also showed us that we are able to perceive the temporal evolution of decorrelated audio features simultaneously. This multidimensionality of timbre perception has been extensively demonstrated in psychoacoustic studies [154] through the concept of multidimensional timbre spaces. Several authors had already pointed out that timbre is a multidimensional phenomenon [304, 269], and that our perception organizes these dimensions given a sound context [400]. Most of the psychoacoustic studies that can be found in the literature often makes use of these multidimensional spaces in order to distinguish different sounds [339, 400]. Timbre has even been said to “*tend to be the psychoacoustician’s multidimensional wastebasket category for everything that cannot be labeled pitch or loudness*” [263]. Authors in the Music Information Retrieval (MIR) community have also pointed out the multifaceted nature of audio perception [114] and that a single measure is unlikely to convey the perceptual similarity of audio signals [386]. Sound retrieval systems should thus be flexible enough so that depending on listeners and target timbres, variable importance could be put on different sound properties during perceptual similarity evaluations [264], but yet no current audio-retrieval system seems to address these limitations.

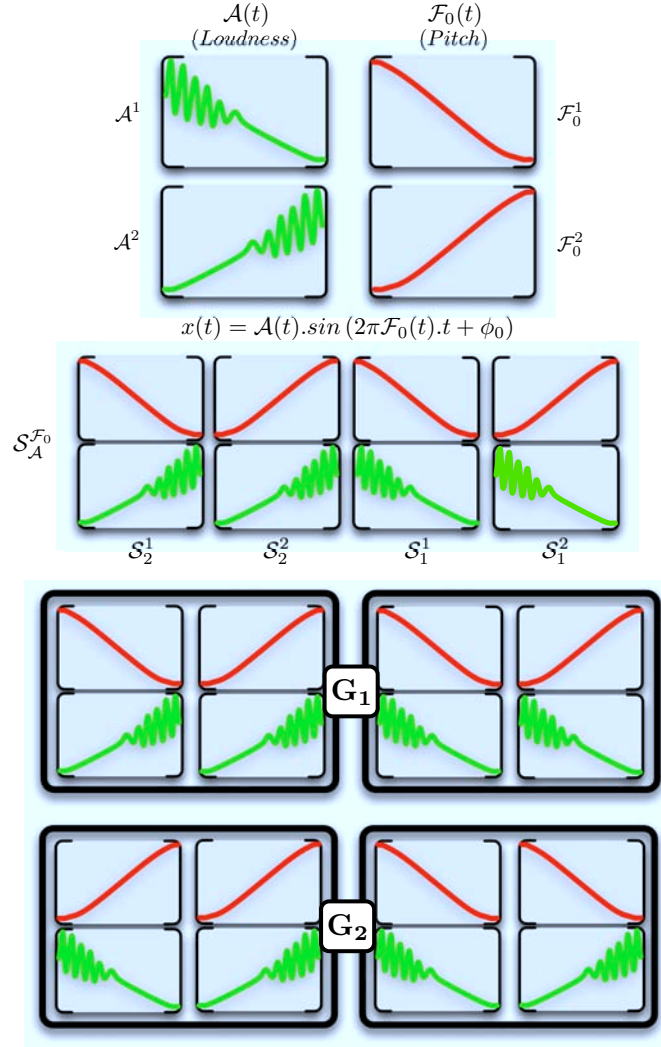


Figure 5: The assessment of sound similarity given its most elementary unit, a synthesized sinusoidal signal. For this signal, we can define the temporal function of its *amplitude* ("loudness") and its *frequency* ("pitch"). Therefore, we define two very simple temporal functions for each feature. The loudness $A(t)$ can either be set to A^1 or A^2 and the pitch $F_0(t)$ can either be set to F_0^1 or F_0^2 . Given these two possibilities for each feature, we can easily synthesize the set $S_A^{F_0}$ of four different sounds S_1^1 , S_2^1 , S_1^2 and S_2^2 . If we try to assess the similarity between elements inside this set, it seems unfeasible to choose if we should to group elements based on the evolution of their *loudness* properties such as grouping G_1 or rather based on the temporal evolution of their *pitch* such as grouping G_2 .

4

PUTTING THE PIECES TOGETHER

4.1 RATIONALE OF THIS STUDY

We discussed along the previous chapters the artistic problematic that raised the questions of this study. Musical orchestration is the art of combining the timbre of instruments in order to achieve a musical thought. While discussing the intrinsic properties of what we called *musical atoms* to be combined, we saw that we should especially take care of two concerns. First, the temporal evolution of audio properties should be a central element of this study. Second, we discussed the extent of our auditory perception which clearly shows that we are able to discern several decorrelated temporal properties simultaneously. Based on these observations, we anticipate a compelling question that is yet to be answered. In order to obtain a true assessment of the timbral similarity we need to take into account both the multidimensional aspect of timbre and the temporal evolution of its structure. However, to the best of our knowledge, it seems that such an approach has never been undertaken. Therefore, this study focus on providing a flexible framework that can assess the *temporal* similarity between elements on several dimensions, without merging the similarities into a single distance measure.

4.2 EPISTEMOLOGICAL CONSIDERATIONS

We strive to give here a discussion on the epistemological considerations underlying this study. Trying to tackle a complex artistic problematic raised a series of perception issues. From these questions, appeared the need for a paradigm yet to be solved. We will, therefore, focus on trying to formalize this unsolved problem and will present algorithms to solve it. We will see that avoiding to merge the distance measures between the temporal features of each dimension can provide a compelling framework of study. Thus, we will show that the applications of this model can go far beyond the realm of music and audio processing. By applying these concepts to generic classification problems, we will show that these auditory ideas can provide an improvement on a variety of scientific fields. The path that we follow in this study, show us that the artistic problematics can lead to raise questions on a wider scale. Trying to answer these questions can in turn provide comprehensive and powerful approaches that can be beneficial to a variety of scientific fields. Hence, music can help us to study science through the complexity of the problems it contains. We will show that we can also close this epistemological loop, by using all the knowledge gained from these studies in order to solve the first problematic. Figure 6 present this conceptual loop of study.

4.3 SCIENTIFIC CONTRIBUTIONS

We explore in this thesis a wide variety of scientific topics, therefore, we try here to summarize our contributions along this study.

- While reviewing the field of time series data mining, we propose four axioms of robustness through which the robustness of any time series distance measure can be evaluated.

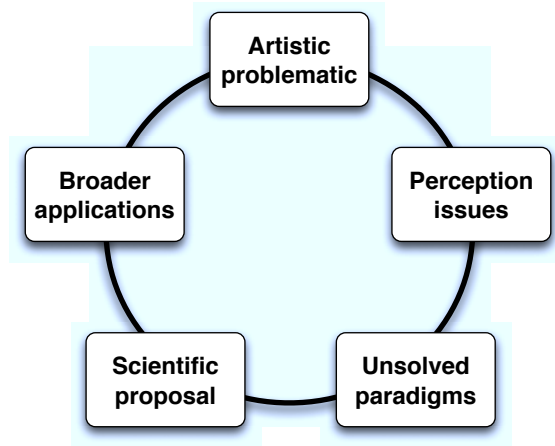


Figure 6: The *epistemological loop* of this study.

- We introduce and formalize the novel problem of *MultiObjective Time Series* (MOTS) matching which allows to take into account both the multidimensional aspects of elements and their temporal structures. We propose two efficient algorithms to solve this problem with a sublinear time complexity. We validate this framework in the field of audio perception through perceptual studies.
- Based on the MOTS framework, we introduce two innovative audio querying paradigm that provide more intuitive interactions with sound samples. We further validate these paradigms with an extensive usability evaluation study.
- We show how to adapt this notion of flexible similarity matching to classification problems. We introduce a new classification selection criteria based on the hyper-volume dominated by each class and study this new classification paradigm over a wide range of datasets. By not merging distances into a single similarity measure, we show that our classification paradigm outperforms state-of-art methods on several scientific fields.
- Based on this paradigm and by still drawing inspiration from the art of listening, we show how to construct a system for biometric identification based on the sounds produced by heart beats. In order to do so, we propose a novel set of features based on the Stockwell transform, called *S-Features*.
- We also show some specific audio applications of our classification paradigm where it allows to outperform state-of-art methods on reference audio classification problems.
- We construct a novel computer-aided orchestration algorithm that can take into account the temporal evolution of audio features. We introduce a multi-level segmentation method that can take into account the different segments that are contained *within* an element at smaller time scales. We show that altogether these new orchestration procedures strongly outperforms the previous approaches.

4.4 STRUCTURE OF THIS DOCUMENT

ELEMENTS OF TIME We start by giving an in-depth review of the time series data mining field (Chapter 5), through its most representative tasks (Section 5.3). We further divide the relevant literature based on the three main aspects of time series handling, namely *representation* methods (Section 5.4.2), *distance* measures (Section 5.4.3) and *indexing* structure (Section 5.4.4). We then provide a overview of the multiobjective optimization field (Chapter 6) by introducing its core notions (Section 6.1), providing a summary of algorithms classification (Section 6.2) and finally presenting some of its application. (Section 6.3).

MULTIOBJECTIVE TIME SERIES (MOTS) MATCHING We then introduce the problem of *MultiObjective Time Series* (MOTS) matching and its formalization (Chapter 7); we further underline the core differences between this novel problem and multivariate matching (Section 7.2) and also briefly discuss its computational complexity. We introduce two efficient algorithms to solve this problem (Section 7.3) and analyze their relative merits on synthetic and real datasets. We describe the application of this framework for innovative audio querying (Section 7.5) by introducing two new querying problematics in the field of audio samples retrieval, namely the *MultiObjective Spectral Evolution Query* (MOSEQ) (Section 7.6.2) and *Query by Vocal Imitation* (QVI) (Section 7.6.3). We validate the hypotheses on which the MOTS framework is based through perceptual studies (Chapter 8) and then perform a comprehensive usability evaluation of the MOSEQ and QVI querying paradigms (Section 7.6.3).

HYPERVOLUME CLASSIFICATION (HV-MOTS) Based on these results, we extend the range of this study and show how to apply these notions of variable similarity evaluation to classification problems by introducing the HyperVolume-MOTS (HV-MOTS) classification scheme (Chapter 10). We show that even within the multiobjective framework that avoids merging distances into a single measure, we can still rank classes by relying on the hypervolume dominated by each (Section 10.1). We discuss the relationship between this novel classification scheme and other distance-based classifiers (Section 10.2) but also to more generic classification schemes (Section 10.3) and discuss its main advantages and drawbacks. We provide a large scale study of the performances of this classification technique (Chapter 11) on a wide range of datasets that covers several scientific fields. We show the statistical superiority of the HV-MOTS classifier over well-established classification schemes (Section 11.3) and state-of-art results on the same datasets (Section 11.4). Based on the HV-MOTS classifier, we show how to build a biometric identification system for heart beat sounds (Chapter 12). We construct it by considering listening as an art (Section 12.3) and developing a specific set of features based on the Stockwell transform, called *S-Features* (Section 12.3.3). We show that using heart sounds as a biometric feature provide a reliable identification (Section 12.4.3) and that this feature is not affected by the phenomenon of *template ageing* (Section 12.4.3) over a time span of two years, supported by the recordings collected in the Mars 500 isolation study. We illustrate the application of the HV-MOTS framework to audio problems (Chapter 16) through generic audio samples classification (Section 13.1) and sound morphology (Section 13.2).

GOING BACK TO MUSIC We show how this knowledge gained through broader applications can be put to use in the field of musical orchestration (Chapter 15). We propose a new orchestration system based on an algorithm that use the MOTS

framework and that rely on an entropic segmentation method (Section 15.4). We show that this new algorithm outperforms the previous approaches for computer-aided orchestration (Section 15.4.2). We then present other artistic applications of the MOTS framework (Chapter 16).

Finally, we offer our lines of future work (Chapter 18) and conclusions (Chapter 19).

Part II

ELEMENTS OF TIME AND PERCEPTION

5

TIME SERIES DATA MINING

In almost every scientific field, measurements are performed over time. These observations lead to a collection of organized data called *time series*. The purpose of time series data mining is to try to extract all meaningful knowledge from the *shape* of data. Even if humans have a natural capacity to perform these tasks, it remains a complex problem for computers. In this chapter, we intend to provide a survey of the techniques applied for time series data mining. The first part is devoted to an overview of the tasks that have captured most of the interest of researchers. Considering that in most cases, time series task relies on the same components for implementation, we divide the literature depending on these common aspects, namely *representation* techniques, *distance* measures and *indexing* methods. The study of the relevant literature has been categorized for each individual aspects. We also introduce in this chapter four types of robustness that we formalize and thanks to which any kind of distance measure could then be classified. Finally, we submit various research trends and avenues that can be explored in the near future.

5.1 INTRODUCTION

A time series represents a collection of values obtained from sequential measurements over time. Time series data mining stems from the desire to reify our natural ability to visualize the *shape* of data. Humans rely on complex schemes in order to perform such tasks. We can actually avoid focusing on small fluctuations in order to derive a notion of *shape* and identify almost instantly similarities between patterns on various time scales. Major time series related tasks include query by content [124], anomaly detection [399], motif discovery [250], prediction [398], clustering [247], classification [22] and segmentation [214]. Despite the vast body of work devoted to this topic in the early years, [11] noted that “*the research has not been driven so much by actual problems but by an interest in proposing new approaches*”. However, with the ever-growing maturity of time series data mining techniques, this statement seems to have become obsolete. Nowadays, time series analysis covers a wide range of real-life problems in various fields of research. Some examples include economic forecasting [357], intrusion detection [432], gene expression analysis [252], medical surveillance [64] and hydrology [283].

Time series data mining unveils numerous facets of complexity. The most prominent problems arise from the high dimensionality of time series data and the difficulty of defining a form of similarity measure based on human perception. With the rapid growth of digital sources of information, time series mining algorithms will have to match increasingly massive datasets. These constraints show us that three major issues are involved:

- *Data representation*: How can the fundamental *shape characteristics* of a time series be represented? What invariance properties should the representation satisfy? A representation technique should derive the notion of shape by reducing the dimensionality of data while retaining its essential characteristics.

- *Similarity measurement*: How can any pair of time series be distinguished or matched? How can an intuitive distance between two series be formalized? This measure should establish a notion of similarity based on perceptual criteria, thus allowing the recognition of perceptually similar objects even though they are not mathematically identical.
- *Indexing method*: How should a massive set of time series be organized to enable fast querying? In other words, what *indexing mechanism* should be applied? The indexing technique should provide minimal space consumption and computational complexity.

These implementation components represent the core aspects of time series data mining systems. However these are not exhaustive as many tasks will require the use of more specific modules. Moreover, some of these are useless for some specific tasks. Forecasting (cf. section 5.3.5) is the most blatant example of a topic that requires more advanced analysis processes as it is more closely related to statistical analysis. It may require the use of a time series representation and a notion of similarity (mostly used to measure prediction accuracy) whereas model selection and statistical learning are also at the core of forecasting systems. The components that are *common* to most time series mining tasks have therefore been analyzed and other components found in related tasks have been briefly discussed.

5.2 DEFINITIONS

The purpose of this section is to provide a definition for the terms used throughout our study.

Definition 1. A *time series* T is an ordered sequence of n real-valued variables

$$T = (t_1, \dots, t_n), t_i \in \mathbb{R} \quad (5.1)$$

A time series is often the result of the observation of an underlying process in the course of which values are collected from measurements made at uniformly spaced *time instants* and according to a given *sampling rate*. A time series can thus be defined as a set of contiguous time instants. The series can be *univariate* as in definition 1 or *multivariate* when several series simultaneously span multiple dimensions within the same time range.

Time series can cover the full set of data provided by the observation of a process and may be of considerable length. In the case of streaming, they are semi-infinite as time instants continuously feed the series. It thus becomes interesting to consider only the *subsequences* of a series.

Definition 2. Given a time series $T = (t_1, \dots, t_n)$ of length n , a *subsequence* S of T is a series of length $m \leq n$ consisting of contiguous time instants from T

$$S = (t_k, t_{k+1}, \dots, t_{k+m-1}) \quad (5.2)$$

with $1 \leq k \leq n - m + 1$. We denote the set of all subsequences of length m from T as S_T^m .

For easier storage, massive time series sets are usually organized in a database.

Definition 3. A *time series database* DB is an unordered set of time series.

As one of the major issues with time series data mining is the *high dimensionality* of data, the database usually contains only simplified representations of the series.

Definition 4. Given a time series $T = (t_1, \dots, t_n)$ of length n , a *representation* of T is a model \bar{T} of reduced dimensionality \bar{d} ($\bar{d} \ll n$) such that \bar{T} closely approximates T .

Nearly every task of time series data mining relies on a notion of similarity between series. After defining the general principle of similarity measures between time series, we will see (section 5.4.3) how these can be specified.

Definition 5. The *similarity measure* $\mathcal{D}(T, U)$ between time series T and U is a function taking two time series as inputs and returning the *distance* d between these series.

This distance has to be *non-negative*, i.e. $\mathcal{D}(T, U) \geq 0$. If this measure satisfies the additional *symmetry* property $\mathcal{D}(T, U) = \mathcal{D}(U, T)$ and *subadditivity* $\mathcal{D}(T, V) \leq \mathcal{D}(T, U) + \mathcal{D}(U, V)$ (also known as the *triangle inequality*), the distance is said to be a *metric*. As will be seen below (section 5.4.4), on the basis of the triangle inequality, metrics are very efficient measures for indexing. We may also extend this notion of distance to the subsequences.

Definition 6. The *subsequence similarity measure* $\mathcal{D}_{\text{subseq}}(T, S)$ is defined as

$$\mathcal{D}_{\text{subseq}}(T, S) = \min(\mathcal{D}(T, S')) \quad (5.3)$$

for $S' \in \mathbf{S}_S^{|T|}$. It represents the distance between T and its best matching location in S .

5.3 TASKS IN TIME SERIES DATA MINING

This section provides an overview of the tasks that have attracted wide research interest in time series data mining. These tasks are usually just defined as theoretical objectives though concrete applications may call for simultaneous use of multiple tasks.

5.3.1 Query by content

Query by content is the most active area of research in time series analysis. It is based on retrieving a set of solutions that are most similar to a query provided by the user. Figure 7 depicts a typical query by content task, represented on a 2-dimensional search space. We can define it formally as

Definition 7. *Query by content* - Given a query time series $Q = (q_1, \dots, q_n)$ and a similarity measure $\mathcal{D}(Q, T)$, find the ordered list $\mathcal{L} = \{T_1, \dots, T_n\}$ of time series in the database DB , such that $\forall T_k, T_j \in \mathcal{L}, k > j \Rightarrow \mathcal{D}(Q, T_k) > \mathcal{D}(Q, T_j)$.

The content of the result set depends on the *type* of query performed over the database. The previous definition is in fact a generalized formalization of a query by content. It is possible to specify a threshold ϵ and retrieve all series whose similarity with the query $\mathcal{D}(Q, T)$ is less than ϵ . This type of query is called an *ϵ -range query*.

Definition 8. *ϵ -range query* - Given a query time series $Q = (q_1, \dots, q_n)$, a time series database DB , a similarity measure $\mathcal{D}(Q, T)$ and a threshold ϵ , find the set of series $\mathcal{S} = \{T_i \mid T_i \in DB\}$ that are within distance ϵ from Q . More precisely, find $\mathcal{S} = \{T_i \in DB \mid \mathcal{D}(Q, T_i) \leq \epsilon\}$

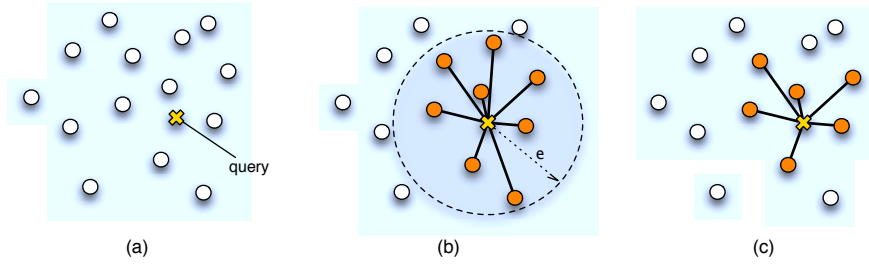


Figure 7: Diagram of a typical query by content task represented in a 2-dimensional search space. Each point in this space represents a series whose coordinates are associated with its features. (a) When a query is entered into the system, it is first transformed into the same representation as that used for other datapoints. Two types of query can then be computed. (b) A ϵ -range query will return the set of series that are within distance ϵ of the query. (c) A K -Nearest Neighbors query will return the K points closest to the query.

Selecting this threshold is obviously highly data-dependent. Users usually want to retrieve a set of solutions by constraining the number of series it should contain, without knowing how far they will be from the query. It is thus possible to query the K most similar series in the database (K -Nearest Neighbors query).

Definition 9. *K-Nearest Neighbors* - Given a query time series $Q = (q_1, \dots, q_n)$, a time series database DB , a similarity measure $\mathcal{D}(Q, T)$ and an integer K , find the set of K series that are the most similar to Q . More precisely, find $\mathcal{S} = \{T_i \mid T_i \in DB\}$ such that $|\mathcal{S}| = K$ and $\forall T_j \notin \mathcal{S}, \mathcal{D}(Q, T_i) \leq \mathcal{D}(Q, T_j)$

Such queries can be called on complete time series; however, the user may also be interested in finding every subsequence of the series matching the query, thus making a distinction between *whole series matching* and *subsequence matching*.

Definition 10. *Whole series matching* - Given a query Q , a similarity measure $\mathcal{D}(Q, T)$ and a time series database DB , find all series $T_i \in DB$ such that $\mathcal{D}(Q, T_i) \leq \epsilon$

This distinction between these types of queries is here expressed in terms of ϵ -range query

Definition 11. *Subsequence matching* - Given a query Q , a similarity measure $\mathcal{D}(Q, T)$ and a database DB , find all subsequences T'_i of series $T_i \in DB$ such that $\mathcal{D}_{\text{subseq}}(Q, T'_i) \leq \epsilon$

In former times, time series mining was almost exclusively devoted to this task (cf. seminal work by [5]). In this paper, the representation was based on a set of coefficients obtained from a Discrete Fourier Transform (DFT) to reduce the dimensionality of data. These coefficients were then indexed with a R^* -tree [31]. False hits were removed in a post-processing step, applying the Euclidean distance to complete time series. This paper laid the foundations of a reference framework that many subsequent works just enlarged by using properties of the DFT [316] or similar decompositions such as Discrete Wavelet Transform (DWT) [83], that has been shown to have similar efficiency depending on the dataset at hand [305]. The Discrete Cosine Transform (DCT) has also been suggested [226] but it appeared later that it did not have any advantage over other decompositions [216]. Several numeric transformations – such as random projections [186], Piecewise Linear Approximation (PLA) [350], Piecewise Approximate Aggregation (PAA) [213, 418] and Adaptive Piecewise Constant Approximation (APCA)

[212] – have been used as representations. Symbolic representations have also been widely used. A shape alphabet with fixed resolution was originally proposed in [6]. Other symbolic representations have been proposed, such as the bit level approximation [320] or the Symbolic Aggregate approXimation (SAX) [249]; the latter one has been shown to outperform most of the other representations [360]. We will find below a detailed overview of representations (section 5.4.2), distance measures (section 5.4.3) and indexing techniques (section 5.4.4).

Other important extensions to query by content include the handling of scaling and gaps [387], noise [389], query constraints [147] and time warping, either by allowing false dismissals [419] or working without constraints [334]. Lower bounding distances without false dismissals for DTW were proposed in [220] and [211] which allows exact indexing. The recent trend of query by content systems seems to be focused on streams. Given the continuously growing bandwidth, most of next generation analysis will most likely have to be performed over stream data. The dynamic nature of streaming time series precludes using the methods proposed for the static case. In a recent study, [225] introduced the most important issues concerning similarity search in static and streaming time series databases. In [224], the use of an incremental computation of DFT allows to adapt to the stream update frequency. However, maintaining the indexing tree for the whole streaming series seems to be uselessly costly. [14] proposed a filter-and-refine DTW algorithm called Anticipatory DTW, which allows faster rejection of false candidates. [244] proposed a weighted locality-sensitive hashing (WLSH) technique applying to approximate queries and working by incremental updating adaptive to the characteristics of stream data. [243] proposed three approaches, polynomial, DFT and probabilistic, to predict future unknown values and answer queries based on the predicated data. This approach is a combination of prediction (cf. section 5.3.5) and streaming query by content; it is representative of an effort to obtain a convergence of approaches that seem to be heterogeneous.

5.3.2 Clustering

Clustering is the process of finding natural groups, called *clusters*, in a dataset. The objective is to find the most homogeneous clusters that are as distinct as possible from other clusters. More formally, the grouping should maximize inter-cluster variance while minimizing intra-cluster variance. The algorithm should thus automatically locate which groups are intrinsically present in the data. Figure 8 depicts some possible outputs of a clustering algorithm. It can be seen in this figure that the main difficulty concerning any clustering problem (even out of the scope of time series mining) usually lies in defining the correct number of clusters. The time series clustering task can be divided into two sub-tasks.

Whole series clustering

Clustering can be applied to each complete time series in a set. The goal is thus to regroup entire time series into clusters so that the time series are as similar to each other as possible within each cluster.

Definition 12. Given a time series database DB and a similarity measure $\mathcal{D}(Q, T)$, find the set of clusters $\mathcal{C} = \{c_i\}$ where $c_i = \{T_k \mid T_k \in \text{DB}\}$ that maximizes inter-cluster distance and minimizes intra-cluster variance. More formally $\forall i_1, i_2, j$ such that $T_{i_1}, T_{i_2} \in c_i$ and $T_j \in c_j$ $\mathcal{D}(T_{i_1}, T_j) \gg \mathcal{D}(T_{i_1}, T_{i_2})$

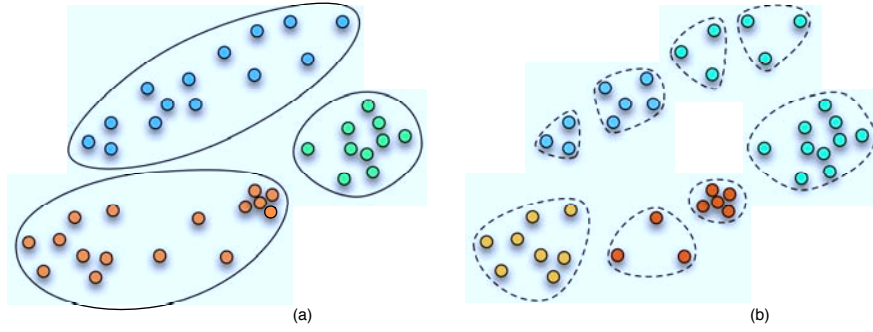


Figure 8: Two possible outputs from the same clustering system obtained by changing the required number of clusters with (a) $N = 3$ and (b) $N = 8$. As we can see, the clustering task is a non trivial problem that highly depends on the way parameters are initialized and the level of detail targeted. This parameter selection issue is common to every clustering task, even out of the scope of time series mining.

There have been numerous approaches for whole series clustering. Typically, after defining an adequate distance function, it is possible to adapt any algorithm provided by the generic clustering topic. Clustering is traditionally performed by using Self Organizing Maps (SOM) [86], Hidden Markov Models (HMM) [356] or Support Vector Machines (SVM) [420]. [140] proposed a variation of the Expectation Maximization (EM) algorithm. However, this model-based approach has usually some scalability problems and implicitly presupposes the existence of an underlying model which is not straightforward for every dataset. Using Markov chain Monte Carlo (MCMC) methods, [135] makes an estimation about the appropriate grouping of time series simultaneously along with the group-specific model parameters. A good survey of generic clustering algorithms from a data mining perspective is given in [39]. This review focuses on methods based on classical techniques that can further be applied to time series. A classification of clustering methods for various static data is proposed in [166] following five categories: *partitioning*, *hierarchical*, *density-based*, *grid-based* and *model-based*. For the specificities of time series data, three of these five categories (partitioning, hierarchical and model-based) have been applied [245]. Clustering of time series is especially useful for data streams; it has been implemented by using clipped data representations [18], Auto-Regressive (AR) models [98], k-Means [388] and – with greater efficiency – k-center clustering [99]. Interested readers may refer to [245] who provides a thorough survey of time series clustering issues by discussing the advantages and limitations of existing works as well as avenues for research and applications.

Subsequence clustering

In this approach, the clusters are created by extracting subsequences from a single or multiple longer time series.

Definition 13. Given a time series $T = (t_1, \dots, t_n)$ and a similarity measure $\mathcal{D}(Q, C)$, find the set of clusters $\mathcal{C} = \{c_i\}$ where $c_i = \{T_j' \mid T_j' \in S_T^n\}$ is a set of subsequences that maximizes inter-cluster distance and intra-cluster cohesion.

In [170], the series are sliced into non-overlapping windows. Their width is chosen by investigating the periodical structure of the time series by means of a DFT analysis. This approach is limited by the fact that, when no strong periodical structure is present in the

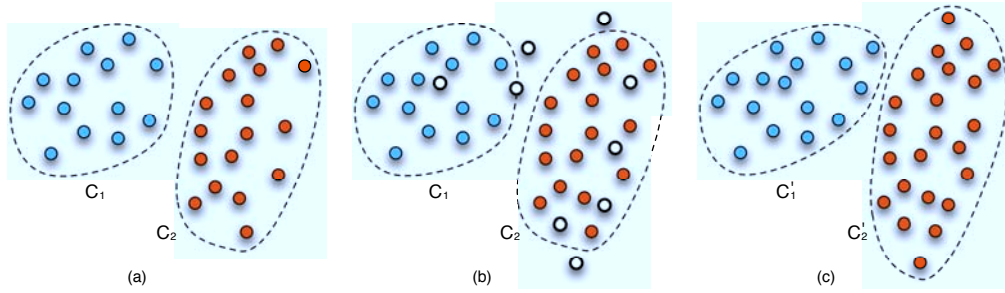


Figure 9: The three main steps of a classification task. (a) A training set consisting of two pre-labeled classes C_1 and C_2 is entered into the system. The algorithm will first try to learn what the characteristic features distinguishing one class from another are; they are represented here by the class boundaries. (b) An unlabeled dataset is entered into the system that will then try to automatically deduce which class each datapoint belongs to. (c) Each point in the set entered has been assigned to a class. The system can then optionally adapt the classes boundaries.

series, non-overlapping slicing may miss important structures. A straightforward way to extend this approach can therefore be to extract shorter overlapping subsequences and then cluster the resulting set. However, this overlapping approach has been shown to produce meaningless results [215]. Despite these deceptive results, the authors pointed out that a meaningful subsequence clustering system could be constructed on top of a motif mining [293] algorithm (cf. section 5.3.7). [107] was first to suggest an approach to overcome this inconsistency by not forcing the algorithm to use all subsequences in the clustering process. In the context of intrusion detection, [432] studied multiple centroid-based unsupervised clustering algorithms, and proposed a self-labeling heuristic to detect any attack within network traffic data. Clustering is also one of the major challenges in bioinformatics, especially in DNA analysis. [218] surveyed state-of-the-art applications of gene expression clustering and provided a framework for the evaluation of results.

5.3.3 Classification

The classification task seeks to assign labels to each series of a set. The main difference when compared to the clustering task is that classes are known in advance and the algorithm is trained on an example dataset. The goal is first to learn what the distinctive *features* distinguishing classes from each others are. Then, when an unlabeled dataset is entered into the system, it can automatically determine which class each series belongs to. Figure 9 depicts the main steps of a classification task.

Definition 14. Given an unlabeled time series T , assign it to one class c_i from a set $\mathcal{C} = \{c_i\}$ of predefined classes.

There are two types of classification. The first one is the *time series classification* similar to whole series clustering. Given sets of time series with a label for each set, the task consists in training a classifier and labeling new time series. An early approach to time series classification was presented in [22]. However, it is based on simple trends whose results are therefore hard to interpret. A piecewise representation was later proposed in [210], it is robust to noise and weighting can be applied in a relevance feedback framework. The same representation was used in [142]; it is apparently not

too robust to outliers. To overcome the obstacle of high dimensionality, [197] used Singular Value Decomposition to select essential frequencies. However, it implies higher computational costs. In a recent study, [330] compared three types of classifiers: nearest neighbor, support vector machines and decision forests. All three methods seems to be valid, though highly depending on the dataset at hand. 1-NN classification algorithm with DTW seems to be the most widely used classifier; it was shown to be highly accurate [405], though computing speed is significantly affected by repeated DTW computations. To overcome this limitation [359] proposed a template construction algorithm based on the Accurate Shape Averaging (ASA) technique. Each training class is represented by only one sequence so that any incoming series is compared only with one averaged template per class. Several other techniques have been introduced, such as ARMA models [106] or HMM [431]. In the context of clinical studies, [252] enhanced HMM approaches by using discriminative HMMs in order to maximize inter-classes differences. Using the probabilistic transitions between fewer states results in the patients being aligned to the model and can account for varying rates of progress. This approach has been applied in [255], in order to detect post-myocardial infarct patients. Several machine learning techniques have also been introduced such as neural networks [280] or Bayesian classification [306]. However, many of these proposals have been shown to be overpowered by a simple 1NN-DTW classifier [405]. A double-loop EM algorithm with a Mixture of Experts network structure has been introduced in [363] for the detection of epileptic seizure based on the EEG signals displayed by normal and epileptic patients. A well-known problem in classification tasks is the *overtraining*, i.e. when too many training data lead to an over-specified and inefficient model. [319] suggested a stopping criterion to improve the data selection during a self training phase. [428] proposed a time series reduction, which extracts patterns that can be used as inputs to classical machine-learning algorithms. Many interesting applications to this problem have been investigated such as brain-computer interface based on EEG signals; they have been reviewed in [254].

5.3.4 Segmentation

The segmentation (or *summarization*) task aims at creating an accurate approximation of time series, by reducing its dimensionality while retaining its essential features. Figure 10 shows the output of a segmentation system. Section 5.4.2 will show that most time series representations try to solve this problem implicitly.

Definition 15. Given a time series $T = (t_1, \dots, t_n)$, construct a model \bar{T} of reduced dimensionality \bar{d} ($\bar{d} \ll n$) such that \bar{T} closely approximates T . More formally $|R(\bar{T}) - T| < \epsilon_r$, $R(\bar{T})$ being the reconstruction function and ϵ_r an error threshold.

The objective of this task is thus to minimize the reconstruction error between a reduced representation and the original time series. The main approach that have been undertaken over the years seems to be Piecewise Linear Approximation (PLA) [350]. The main idea behind PLA is to split the series into most representative segments, and then fit a polynomial model for each segment. A good review on the most common segmentation methods in the context of PLA representation can be found in [214]. Three basic approaches are distinguished. In *sliding windows*, a segment is grown until it exceeds some error threshold [350]. This approach has shown poor performance with many real life datasets [214]. The *top-down* approach consists in recursively partitioning a time series until some stopping criterion is met [239]. This approach has time complexity $O(n^2)$ [291] and is qualitatively outperformed by *bottom-up*. In this

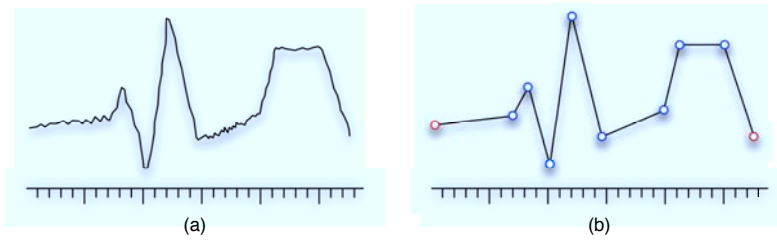


Figure 10: Example of application of a segmentation system. From (a) usually noisy time series containing a very large number of datapoints, the goal is to find (b) the closest approximation of the input time series with the maximal dimensionality reduction factor without losing any of its essential features.

approach, starting from the finest approximation, segments are iteratively merged [210]. [177] present fast greedy algorithms to improve previous approaches and a statistical method for choosing the number of segments is described in [385].

Several other methods have been introduced to handle this task. [286] introduced a representation of time series that implicitly handles the segmentation of time series. They proposed user-specified amnesic functions reducing the confidence to older data in order to make room for newer data. In the context of segmenting hydrological time series, [207] proposed a maximum likelihood method using an HMM algorithm. However, this method offers no guarantee to yield the globally optimal segmentation without long execution times. For dynamic summary generation, [281] proposed an online transform-based summarization techniques over data streams that can be updated continuously. The segmentation of time-series can also be seen as a constrained clustering problem. [2] proposed to group time points by their similarity, provided that all points in a cluster come from contiguous time instants. Therefore, each cluster represents the segments in time whose homogeneity is evaluated with a local PCA model.

5.3.5 Prediction

Time series are usually very long and considered *smooth*, i.e. subsequent values are within predictable ranges of one another [349]. The task of prediction is aimed at explicitly modeling such variable dependencies to forecast the next few values of a series. Figure 11 depicts various forecasting scenarios.

Definition 16. Given a time series $T = (t_1, \dots, t_n)$, predict the k next values $(t_{n+1}, \dots, t_{n+k})$ that are most likely to occur.

Prediction is a major area in several fields of research. Concerning time series, it is one of the most extensively applied tasks. Literature about this is so abundant that dozens of reviews can focus on only a specific field of application or family of learning methods. Even if it can use time series representations and a notion of similarity to evaluate accuracy, It also relies on several statistical components that are out of the scope of this article, e.g. model selection and statistical learning. This task will be mentioned because of its importance but the interested reader willing to have further information may consult several references on forecasting [59, 167, 379, 60] Several methods have been applied to this task. A natural option could be AR models [56]. These models have been applied for a long time to prediction tasks involving signal de-noising or dynamic systems modeling. It is however possible to use more

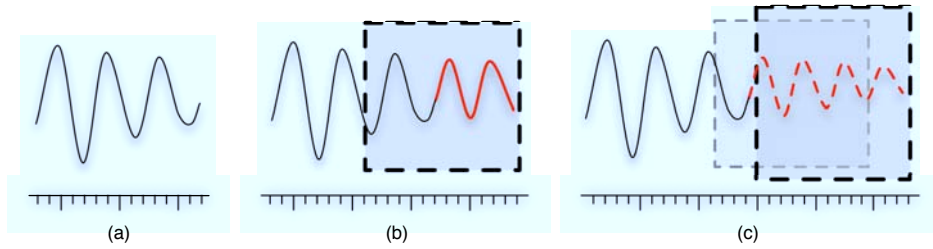


Figure 11: A typical example of the time series prediction task. (a) The input time series may exhibit a periodical and thus predictable structure. (b) The goal is to forecast a maximum number of upcoming datapoints within a prediction window. (c) The task becomes really hard when it comes to having *recursive prediction*, i.e. the long term prediction of a time series implies reusing the earlier forecast values as inputs in order to go on predicting.

complex approaches such as neural networks [227] or clusters function approximation [347] to solve this problem. A polynomial architecture has been developed to improve a multilayer neural network in [411] by reducing higher-order terms to a simple product of linear functions. Other learning algorithms, such as SOM, provided efficient supervised architectures. A survey of applications of SOM to time series prediction is given in [27]. Recent improvements for time series forecasting have been proposed; [300] proposed a Bayesian prediction for time series subject to discrete breaks, handling the size and duration of possible breaks by means of a hierarchical HMM. A dynamic genetic programming (GP) model tailored for forecasting streams was proposed in [393] by adapting incrementally based on retained knowledge. The prediction task seems one of the most commonly applied in real-life applications, considering that market behavior forecasting relies on a wealth of financial data. [21] proposed to refine the method of factor forecasting by introducing ‘targeted predictors’ selected by using a hysteresis (hard and soft thresholding) mechanism. The prediction task has also a wide scope of applications ranging from tourism demand forecasting [357] to medical surveillance [64]. In this paper, the authors compared the predictive accuracy of three methods, namely, non-adaptive regression, adaptive regression, and the Holt-Winters method; the latter appeared to be the best method. In a recent study, [7] carried out a large scale comparison for the major machine-learning models applied to time series forecasting, following which the best two methods turned out to be multilayer perceptron and Gaussian process regression. However, learning a model for long-term prediction seems to be more complicated, as it can use its own outputs as future inputs (*recursive prediction*). [176] proposed the use of least-squares SVM, to solve this problem. [70] further applied saliency analysis to SVM in order to remove irrelevant features based on the sensitivity of the network output to the derivative of the feature input. [358] proposed to combine direct prediction and an input selection in order to cope with long-term prediction of time series.

5.3.6 Anomaly detection

The detection of anomalies seeks to find abnormal subsequences in a series. Figure 12 depicts an example of anomaly detection. It has numerous applications ranging from biosurveillance [95] to intrusion detection [432].

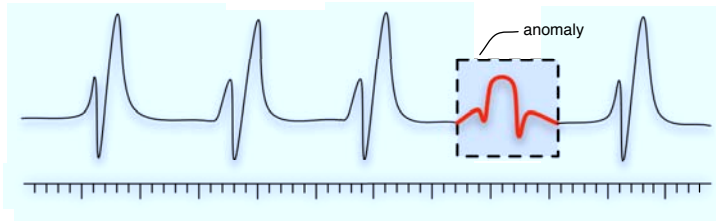


Figure 12: An idealized example of the anomaly detection task. A long time series which exhibits some kind of periodical structure can be modeled thanks to a reduced pattern of “standard” behavior. The goal is thus to find subsequences which does not follow the model and may therefore be considered as anomalies.

Definition 17. Given a time series $T = (t_1, \dots, t_n)$ and a model of its normal behavior, find all subsequences $T' \in S_T^n$ which contain anomalies, i.e. do not fit the model.

A good discussion on the difficulties of mining rare events is given in [399]. The usual approach to detect anomalies is to first create a model of a series’ normal behavior and characterize subsequences that stray too far from the model as anomalies. This approach can be linked to the prediction task. Indeed, if we can forecast the next values of a time series with a large accuracy, outliers can be detected in a straightforward manner and flagged as anomalies. This approach was undertaken first in [422] using SOM model to represent the expected behavior. A framework for novelty detection is defined in [257] and implemented based on Support Vector Regression (SVR). Machine learning techniques were also introduced to dynamically adapt their modelisation of normal behavior. [8] investigated the use of block-based One-Class Neighbor Machine and recursive Kernel-based algorithms and showed their applicability to anomaly detection. [90] proposed two algorithms to find anomalies in the Haar wavelet coefficients of the time series. A state-based approach is taken in [336] using time point clustering so that clusters represents the normal behavior of a series. Another definition of anomalies, the time series *discords*, are defined as subsequences that are maximally different from all the remaining subsequences [217]. This definition is able to capture the idea of most unusual subsequence within a time series and its unique parameter is the required length of the subsequences. Thanks to this definition [414] proposed an exact algorithm that requires only two linear scans, thus allowing for the use of massive datasets. However, as several proposals, the number of anomalous subsequences must be specified prior to the search. Several real-life applications have also been outlined in recent research. Anomaly detection is applied in [162] to detect fatigue damage in polycrystalline alloys, thus preventing problems in mechanical structures. An anomaly detection scheme for time series is used in [95] to determine whether streams coming from sensors contain any abnormal heartbeats. A recent overview and classification of the research on anomaly detection is presented in [84], which provides a discussion on the computational complexity of each technique.

5.3.7 Motif discovery

Motif discovery consists in finding every subsequences (named *motif*) that appears recurrently in a longer time series. This idea was transferred from gene analysis in bioinformatics. Figure 13 depicts a typical example of motif discovery. Motifs were defined originally in [293] as *typical* non-overlapping subsequences. More formally

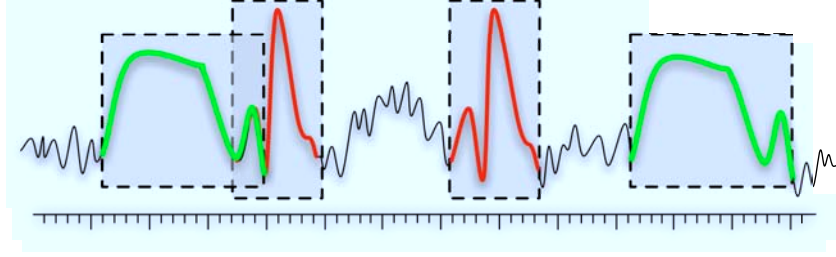


Figure 13: The task of motif discovery consists in finding every subsequence that appears recurrently in a longer time series. These subsequences are named motifs. This task exhibits a high combinatorial complexity as several motifs can exist within a single series, motifs can be of various lengths and even overlap.

Definition 18. Given a time series $T = (t_1, \dots, t_n)$, find all subsequences $T' \in S_T^n$ that occurs repeatedly in the original time series.

A great interest for this research topic has been triggered by the observation that subsequence clustering produces meaningless results [215]. The authors pointed out that motif discovery could be used as a subroutine to find meaningful clusters. In order to find motifs more efficiently, [94] proposed to use the random projection algorithm [63] which was successfully used for DNA sequences. However, because of its probabilistic nature, it is not guaranteed to find the exact set of motifs. [127] proposed an algorithm that can extract approximate motifs in order to mine time series data from protein folding/unfolding simulations. In [253], motif discovery is formalized as a continuous top-k motif balls problem in an m -dimensional space. However, the efficiency of this algorithm critically depends on setting the desired length of the pattern. [372] introduced a k -motif-based algorithm that provides an interesting mechanism to generate summaries of motifs. [413] showed that motif discovery can be severely altered by any slight change of *uniform scaling* (linear stretching of the pattern length) and introduced a scaling-invariant algorithm to determine the motifs. An algorithm for exact discovery of time series motifs has been recently proposed [277], which is able to process very large datasets by using early abandoning on a linear re-ordering of data. [274] studied the constrained motif discovery problem which provides a way to incorporate prior knowledge into the motif discovery process. They showed that most unconstrained motif discovery problems can be transformed into constrained ones and provided two algorithms to solve such problem. The notion of motifs can be applied to many different tasks. The modeling of normal behavior for anomaly detection (cf. section 5.3.6) implies finding the recurrent motif of a series. For time series classification, significant speed-ups can be achieved by constructing motifs for each class [428].

5.4 IMPLEMENTATION COMPONENTS

In this section, we review the implementation components common to most of time series mining tasks. As said earlier, the three key aspects when managing time series data are *representation* methods, *similarity* measures and *indexing* techniques. Because of the *high dimensionality* of time series, it is crucial to design low-dimensional representations that preserve the fundamental characteristics of a series. Given this representation scheme, the *distance* between time series needs to be carefully defined in order to

exhibit perceptually relevant aspects of the underlying similarity. Finally the indexing scheme must allow to efficiently manage and query evergrowing massive datasets.

5.4.1 Preprocessing

In real-life scenarios, time series usually come from live observations [325] or sensors [360] which are particularly subject to noise and outliers. These problems are usually handled by preprocessing the data. Noise filtering can be handled by using traditional signal processing techniques like digital filters or wavelet thresholding. In [178], Independent Component Analysis (ICA) is used to extract the main mode of the series. As will be explained in section 5.4.2, several representations implicitly handle noise as part of the transformation.

The second issue concerns the scaling differences between time series. This problem can be overcome by a linear transformation of the amplitudes [147]. Normalizing to a fixed range [6] or first subtracting the mean (known as *zero mean / unit variance* [212]) may be applied to both time series, however it does not give the optimal match of two series under linear transformations [12]. In [148] the transformation is sought with optional bounds on the amount of scaling and shifting. However, normalization should be handled with care. As noted by [387], normalizing an essentially flat but noisy series to unit variance will completely modify its nature and normalizing small enough subsequences can provoke all series to look the same [247].

Finally, resampling (or *uniform time warping* [284]) can be performed in order to obtain series of the same length [209]. Down-sampling the longer series has been shown to be fast and robust [12].

5.4.2 Representation

As mentioned earlier, time series are essentially high dimensional data. Defining algorithms that work directly on the raw time series would therefore be computationally too expensive. The main motivation of representations is thus to emphasize the essential characteristics of the data in a concise way. Additional benefits gained are efficient storage, speedup of processing as well as implicit noise removal. These basic properties lead to the following requirements for any representation:

- Significant reduction of the data dimensionality
- Emphasis on fundamental shape characteristics on both *local* and *global* scales
- Low computational cost for computing the representation
- Good reconstruction quality from the reduced representation
- Insensitivity to noise or implicit noise handling

Many representation techniques have been investigated, each of them offering different trade-offs between the properties listed above. It is however possible to classify these approaches according to the kind of transformations applied. In order to perform such classification, we follow the taxonomy of [216] by dividing representations into three categories, namely *non data-adaptive*, *data-adaptive* and *model-based*. Figure 14 synthesizes every reviewed representation based on our classification.

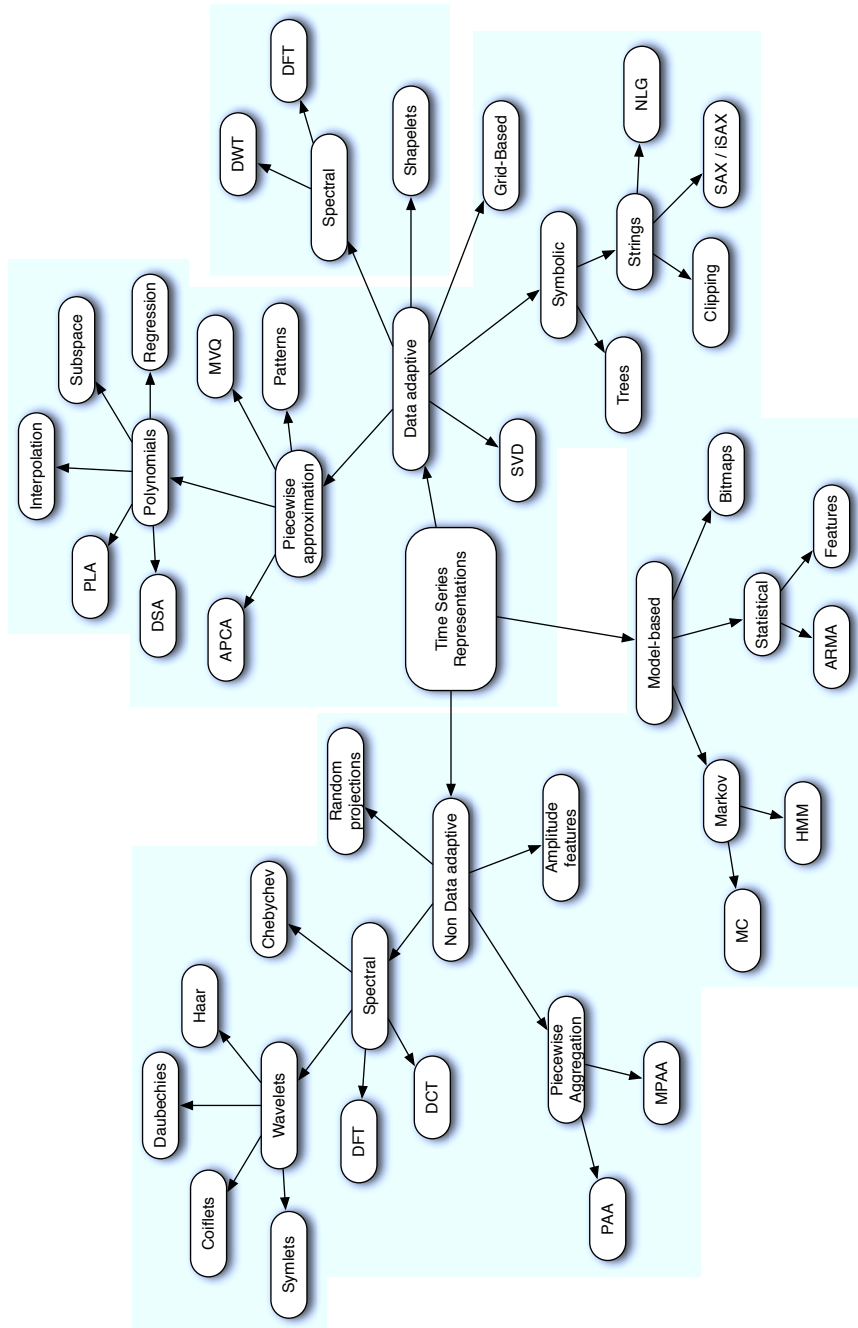


Figure 14: Complete classification of all the time series representations reviewed in this chapter.

Non Data-Adaptive

In non data-adaptive representations, the parameters of the transformation remain the same for every time series regardless of its nature. The first non data-adaptive representations were drawn from spectral decompositions. The DFT was used in the seminal work of [5]. It projects the time series on a sine and cosine functions basis [124] in the real domain. The resulting representation is a set of sinusoidal coefficients. Instead of using a fixed set of basis functions, the DWT uses scaled and shifted versions of a mother *wavelet* function [83]. This gives a multi-resolution decomposition where low frequencies are measured over larger intervals thus providing better accuracy [305]. A large number of wavelet functions have been used in the literature like Haar [82], Daubechies [305] or Coiflets [349]. The Discrete Cosine Transform (DCT) uses only a cosine basis; it has also been applied to time series mining [226]. However, it has been shown that it does not offer any advantage over previously cited decompositions [216]. Finally, an approximation by Chebychev polynomials [66] has also been proposed but the results obtained have later been withdrawn due to an error in implementation.

Other approaches – more specific to time series – have been proposed. The Piecewise Aggregate Approximation (PAA) introduced by [213] (submitted independently as Segmented Means in [418]) represents a series through the mean values of consecutive fixed-length segments. An extension of PAA including a multi-resolution property (MPAA) has been proposed in [247]. [16] suggested to extract a sequence of amplitude-levelwise local features (ALF) to represent the characteristics of local structures. It was shown that this proposal provided weak results in [109]. Random projections have been used for representation in [186]; in this case, each time series enters a convolution product with k random vectors drawn from a multivariate standard. This approach has recently been combined with spectral decompositions by [325] with the purpose of answering statistical queries over streams.

Data-Adaptive

This approach implies that the parameters of a transformation are modified depending on the data available. By adding a data-sensitive selection step, almost all non data-adaptive methods can become data-adaptive. For spectral decompositions, it usually consists in selecting a subset of the coefficients. This approach has been applied to DFT [389] and DWT [362]. A data-adaptive version of PAA has been proposed in [265], with vector quantization being used to create a codebook of recurrent subsequences. This idea has been adapted to allow for multiple resolution levels [266]. However, this approach has only been tested on smaller datasets. A similar approach has been undertaken in [360] with a codebook based on motion vectors being created to spot gestures. However, it has been shown to be computationally less efficient than SAX.

Several inherently data-adaptive representations have also been used. SVD has been proposed [226] and later been enhanced for streams [324]. However, SVD requires computation of eigenvalues for large matrices and is therefore far more expensive than other mentioned schemes. It has recently been adapted to find multi-scale patterns in time series streams [288]. PLA [350] is a widely used approach for the segmentation task (cf. section 5.3.4). The set of polynomial coefficients can be obtained either by interpolation [210] or regression [181]. Many derivatives of this technique have been introduced. The Landmarks system [298] extends this notion to include a multi-resolution property. However, the extraction of features relies on several parameters which are highly data-dependent. APCA [212] uses constant approximations per seg-

ment instead of polynomial fitting. Indexable PLA has been proposed by [89] to speed up the indexing process. [285] put forward an approach based on PLA, to answer queries about the recent past with greater precision than older data and called such representations amnesic. The method consisting in using a segmentation algorithm as a representational tool has been extensively investigated. The underlying idea is that segmenting a time series can be equated with the process of representing the most salient features of a series while considerably reducing its dimensionality. [408] proposed a pattern-based representation of time series. The input series is approximated by a set of concave and convex patterns to improve the subsequence matching process. [425] proposed a pattern representation of time series to extract outlier values and noise. The Derivative Segment Approximation (DSA) model [159] is a representation based on time series segmentation through an estimation of derivatives to which DTW can be applied. The polynomial shape space representation [138] is a subspace representation consisting of trend aspects estimators of a time series. [25] put forward a two-level approach to recognize gestures by describing individual trajectories with key-points, then characterizing gestures through the global properties of the trajectories.

Instead of producing a numeric output, it is also possible to discretize the data into symbols. This conversion into a symbolical representation also offers the advantage of implicitly performing noise removal by complexity reduction. A relational tree representation is used in [23]. Non-terminal nodes of the tree correspond to valleys and terminal nodes to peaks in the time series. The Symbolic Aggregate approXimation (SAX) [249], based on the same underlying idea as PAA, calls on equal frequency histograms on sliding windows to create a sequence of short words. An extension of this approach, called indexable Symbolic Aggregate approXimation (iSAX) [351], has been proposed to make fast indexing possible by providing zero overlap at leaf nodes. The grid-based representation [10] places a two dimensional grid over the time series. The final representation is a bit string describing which values were kept and which bins they were in. Another possibility is to discretize the series to a binary string (a technique called *clipping*) [320]. Each bit indicates whether the series is above or below the average. That way, the series can be very efficiently manipulated. In [20] this is done using the median as the clipping threshold. Clipped series offer the advantage of allowing direct comparison with raw series, thus providing a tighter lower bounding metric. Thanks to a variable run-length encoding, [19] show that it is also possible to define an approximation of the Kolmogorov complexity. Recently, a very interesting approach has been proposed in [417]; it is based on primitives called *shapelets*, i.e. subsequences which are maximally representative of a class and thus fully discriminate classes through the use of a dictionary. This approach can be considered as a step forward towards bridging the gap between time series and shape analysis.

Model-based

The model-based approach is based on the assumption that the time series observed has been produced by an underlying model. The goal is thus to find parameters of such a model as a representation. Two time series are therefore considered similar if they have been produced by the same set of parameters driving the underlying model. Several parametric temporal models may be considered, including statistical modeling by feature extraction [280], ARMA models [202] Markov Chains (MCs) [344] or HMM [287]. MCs are obviously simpler than HMM so they fit well shorter series but their expressive power is far more limited. The Time Series bitmaps introduced in [230] can

also be considered as a model-based representation for time series, even if it mainly aims at providing a visualization of time series.

5.4.3 Similarity measure

Almost every time series mining task requires a subtle notion of similarity between series, based on the more intuitive notion of *shape*. When observing simultaneously multiple characteristics of a series, humans can abstract from such problems as amplitude, scaling, temporal warping, noise and outliers. The Euclidean distance is obviously unable to reach such a level of abstraction. Numerous authors have pointed out several pitfalls when using L_p norms [109, 209, 418]. However, it should be noted that, in the case of very large datasets, Euclidean distance has been shown [351] to be sufficient as there is a larger probability that an almost exact match exists in the database. Otherwise, a similarity measure should be consistent with our intuition and provide the following properties:

1. It should provide a recognition of perceptually similar objects, even though they are not mathematically identical;
2. It should be consistent with human intuition;
3. It should emphasize the most salient features on both *local* and *global* scales;
4. A similarity measure should be universal in the sense that it allows to identify or distinguish arbitrary objects, i.e. no restrictions on time series are assumed;
5. It should abstract from distortions and be invariant to a set of transformations.

Many authors have reported about various transformation invariances required for similarity. Given a time series $T = \{t_1, \dots, t_n\}$ of n datapoints, we consider the following transformations:

- *Amplitude shifting*: The series $G = \{g_1, \dots, g_n\}$ obtained by a linear amplitude shift of the original series $g_i = t_i + k$ with $k \in \mathbb{R}$ a constant.
- *Uniform amplification*: The series G obtained by multiplying the amplitude of the original series $g_i = k \cdot t_i$ with $k \in \mathbb{R}$ a constant.
- *Uniform time scaling*: The series $G = \{g_1, \dots, g_m\}$ produced by a uniform change of the time scale of the original series $g_i = t_{\lceil k \cdot i \rceil}$ with $k \in \mathbb{R}$ a constant.
- *Dynamic amplification*: The series G obtained by multiplying the original series by a dynamic amplification function $g_i = h(i) \cdot t_i$ with $h(i)$ a function such that $\forall t \in [1 \dots n], h'(t) = 0$ if and only if $t'_i = 0$.
- *Dynamic time scaling*: The series G obtained by a dynamic change of the time scale $g_i = t_{h(i)}$ with $h(i)$ a positive and strictly increasing function such $h : \mathbb{N} \rightarrow [1 \dots n]$
- *Additive Noise*: The series G obtained by adding a noisy component to the original series $g_i = t_i + \epsilon_i$ with ϵ_i an independent identically distributed white noise.
- *Outliers*: The series G obtained by adding outliers at random positions. Formally, for a given set of random time positions $\mathcal{P} = \{k \mid k \in [1 \dots n]\}$, $g_k = \epsilon_k$ with ϵ_k an independent identically distributed white noise.

The similarity measure $\mathcal{D}(T, G)$ should be robust to any combinations of these transformations. This property lead to our formalization of four general types of robustness. We introduce properties expressing robustness for *scaling* (amplitude modifications), *warping* (temporal modifications), *noise* and *outliers*. Let \mathcal{S} be a collection of time series, and let \mathcal{H} be the maximal group of homeomorphisms under which \mathcal{S} is closed. A similarity measure \mathcal{D} on \mathcal{S} is called *scale robust* if it satisfies

Proposition. *For each $T \in \mathcal{S}$ and $\alpha > 0$ there is a $\delta > 0$ such that $\|t_i - h(t_i)\| < \delta$ for all $t_i \in T$ implies $\mathcal{D}(T, h(T)) < \alpha$ for all $h \in \mathcal{H}$.*

We call a similarity measure *warp robust* if the following holds

Proposition. *For each $T = \{t_i\} \in \mathcal{S}, T' = \{t_{h(i)}\}$ and $\alpha > 0$ there is a $\delta > 0$ such that $\|i - h(i)\| < \delta$ for all $t_i \in T$ implies that $\mathcal{D}(T, T') < \alpha$ for all $h \in \mathcal{H}$.*

We call a similarity measure *noise robust* if it satisfies the following property

Proposition. *For each $T \in \mathcal{S}$ and $\alpha > 0$, there is a $\delta > 0$ such that $U = T + \epsilon$ with $p(\epsilon) = N(0, \delta)$ implies $\mathcal{D}(T, U) < \alpha$ for all $U \in \mathcal{S}$*

We call a measure *outlier robust* if the following holds

Proposition. *For each $T \in \mathcal{S}, \mathcal{K} = \{\text{rand}[1...n]\}$ and $\alpha > 0$, there is a $\delta > 0$ such that if $|\mathcal{K}| < \delta$ and $U_{k \in \mathcal{K}} = \epsilon_k$ and $U_{k \notin \mathcal{K}} = T_k$ implies $\mathcal{D}(T, U) < \alpha$ for all $U \in \mathcal{S}$*

Similarity measures can be classified in four categories. *Shape-based* distances compare the overall shape of the series. *Edit-based* distances compare two time series on the basis of the minimum number of operations needed to transform one series into another one. *Feature-based* distances extract features describing aspects of the series that are then compared with any kind of distance function. *Structure-based* similarity aims at finding higher-level structures in the series to compare them on a more global scale. We further subdivide this category into two specific subcategories. *Model-based* distances work by fitting a model to the various series and then comparing the parameters of the underlying models. *Compression-based* distances analyze how well two series can be compressed together. Similarity is reflected by higher compression ratios. As defined by [209], we refer to distance measures that compare the i th point of a series to the i th point of another as *lock-step* and measures that allow flexible (one-to-many / one-to-none) comparison as *elastic*. Figure 15 synthesize every reviewed distance measure based on this classification.

Shape-based

The Euclidean distance and other L_p norms [418] have been the most widely used distance measures for time series [209]. However, these have been shown to be poor similarity measurements [11, 109]. As a matter of fact, these measures does not match any of the types of robustness. Even if the problems of scaling and noise can be handled in a preprocessing step [147], the warping and outliers issues need to be addressed with more sophisticated techniques. This is where the use of elastic measures can provide an elegant solution to both problems.

Handling the local distortions of the time axis is usually addressed using *non-uniform time warping* [210], more specifically with Dynamic Time Warping (DTW) [41]. This measure is able to match various sections of a time series by allowing warping of the time axis. The optimal alignment is defined by the shortest warping path in a

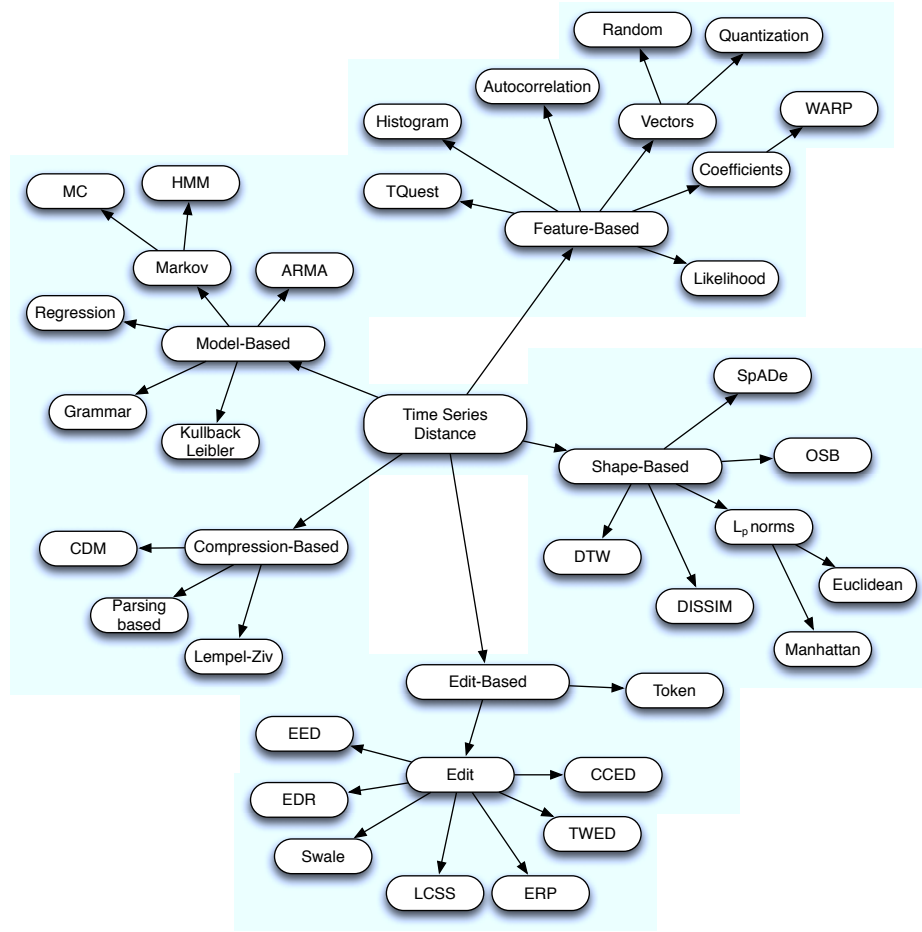


Figure 15: Complete classification of the distance measures reviewed in this chapter.

distance matrix. A warping path W is a set of contiguous matrix indices defining a mapping between two time series. Even if there is an exponential number of possible warping paths, the optimal path is the one that minimizes the global warping cost. DTW can be computed using dynamic programming with time complexity $O(n^2)$ [321]. However, several lower bounding measures have been introduced to speed up the computation. [211] introduced the notion of upper and lower envelope that represents the maximum allowed warping. Using this technique, the complexity becomes $O(n)$. It is also possible to impose a *temporal constraint* on the size of the DTW warping window. It has been shown that these improve not only the speed but also the level of accuracy as it avoids the pathological matching introduced by extended warping [322]. The two most frequently used global constraints are the Sakoe-Chiba Band and the Itakura Parallelogram. [335] introduced the FastDTW algorithm which makes a linear time computation of DTW possible by recursively projecting a warp path to a higher resolution and then refining it. A drawback of this algorithm is that it is approximate and therefore offer no guarantee to finding the optimal solution. In addition to dynamic warping, it may sometimes be useful to allow a global scaling of time series to achieve meaningful results, a technique known as *uniform scaling* (US). [137] proposed the scaled and warped matching (SWM) similarity measure that makes it possible to combine the benefits of DTW with those of US.

Other shape-based measures have been introduced such as the Spatial Assembling Distance (SpADe) [92]; it is a pattern-based similarity measure. This algorithm identifies matching *patterns* by allowing shifting and scaling on both temporal and amplitude axes, thus being scale robust. The DISSIM [134] distance has been introduced to handle similarity at various sampling rates. It is defined as an approximation of the integral of the Euclidean distance. One of the most interesting recent proposals is based on the concept of elastic matching of time series [233]. [234] presented an optimal subsequence matching (OSB) technique that is able to automatically determine the best subsequence and warping factor for distance computation; it includes a penalty when skipping elements. Optimality is achieved through a high computational cost; however, it can be reduced by limiting the skipping range.

Edit-based

Edit-based methods (also known as *Levenshtein distance*) has originally been applied to characterize the difference between two strings. The underlying idea is that the distance between strings may be represented by the minimum number of operations needed to transform one string into another, with insertion, deletion and substitution. The presence of outliers or noisy regions can thus be compensated by allowing gaps in matching two time series. [102] use the Longest Common Subsequence (LCSS) algorithm to tackle this problem. The LCSS distance uses a *threshold parameter* ϵ for point matching and a *warping threshold* δ . A fast approximate algorithm to compute LCSS has been described in [52]. [387] normalized the LCSS similarity by the length of the time series and allowed linear transformations. [391] introduced lower-bounding measure and indexing techniques for LCSS. DTW requires the matched time series to be well aligned and its efficiency deteriorates with noisy data as, when matching all the points, it also matches the outliers distorting the true distance between sequences. LCSS has been shown to be more robust than DTW under noisy conditions [387]; this heavily depends on the threshold setting. [276] proposed the Fast Time Series Evaluation (FTSE) method for computing LCSS. On the basis of this algorithm, they proposed the Sequence Weighted Alignment model (Swale) that extends the ϵ threshold-based scoring techniques to include arbitrary match rewards and gap penalties. The Edit Distance on Real sequence (EDR) [88] is an adaptation of the edit distance to real-valued series. Contrary to LCSS, EDR assign penalties depending on the length of the gaps between the series. The Edit Distance with Real Penalty (ERP) [87] attempts to combine the merits of DTW and edit distance by using a *constant reference point*. For the same purpose, [261] submitted an interesting dynamic programming algorithm called Time Warp Edit Distance (TWED). TWED is slightly different from DTW, LCSS, or ERP algorithms. In particular, it highlights a parameter that controls a kind of stiffness of the elastic measure along the time axis. Another extension to the edit distance has been proposed in [278], it has been called the extended edit distance (EED). Following the observation that the edit distance penalizes all change operations with the same cost, it includes an additional term reflecting whether the operation implied characters that are more frequent, therefore closer in distance. A different approach for constraining the edit operations has been proposed in [93]; it is based on the Constraint Continuous Editing Distance (CCED) that adjusts the potential energy of each sequence to achieve optimal similarity. As CCED does not satisfy triangle inequality, a lower bounding distance is provided for efficient indexing.

Feature-based

These measures rely on the computation of a feature set reflecting various aspects of the series. Features can be selected by using coefficients from DFT [350] or DWT decompositions (cf. section 5.4.2). In [192], a likelihood ratio for DFT coefficients has been shown to outperform Euclidean distance. In [390], a combination of periodogram and autocorrelation functions allows to select the most important periods of a series. This can be extended to carrying out local correlation tracking as proposed in [289].

Concerning symbolic representations, [258] represent each symbol with a random vector and a symbolic sequence by the sum of the vectors weighted by the temporal distance of the symbols. In [129] weighted histograms of consecutive symbols are used as features. The similarity search based on Threshold Queries (TQuEST) [15] use a given threshold parameter τ in order to transform a time series into a sequence of *threshold-crossing* time intervals. It has however been shown to be highly specialized with mitigated results on classical datasets [109]. [29] proposed a Fourier-based approach, called WARP and making the using of the DFT phase possible, this being more accurate for a description of object boundaries.

An approach using ideas from shape and feature-based representations has been described in [266]. Typical local shapes are extracted with vector quantization and the time series are represented by histograms counting the occurrences of these shapes at several resolutions. Multiresolution Vector Quantized (MVQ) approximation keeps both local and global information about the original time series, so that defining a multi-resolution and hierarchical distance function is made possible.

Structure-based

Even if the previously cited approaches have been useful for short time series or subsequences applications, they often fail to produce meaningful results on longer series. This is mostly due to the fact that these distances are usually defined to find *local* similarities between patterns. However, when handling very long time series, it might be more profitable to find similarities on a more *global* scale. Structure-based distances [248] are thus designed to identify higher-level structures in series.

MODEL-BASED Model-based distances offer the additional advantage that prior knowledge about the generating process can be incorporated in the similarity measurement. The similarity can be measured by modeling one time series and determining the likelihood that one series was produced by the underlying model of another. Any type of parametric temporal model may be used. HMM with continuous output values or ARMA models are common choices [409]. However, best results are obtained if the model selected is related to the type of production that generated the data available. In [141], HMMs are combined with a piecewise linear representation. In [287] the distance between HMM is normalized to take into account the quality of fit of the series producing the model. [46] use the similarity-based paradigm where HMM is used to determine the similarity between each object and a pre-determinate set of other objects. For short time series, it is also possible to use regression models as proposed by [140].

Among other common choices for symbolic representations, we may cite MC [327], HMM with discrete output distributions [235], and grammar based models [11]. Alternatively to pairwise likelihood, the Kullback-Leibler divergence allows to have direct comparison of models [344].

COMPRESSION-BASED [216], inspired by results obtained in bioinformatics, defined a distance measure based on the Kolmogorov complexity called Compression-Based Dissimilarity Measure (CDM). The underlying idea is that concatenating and compressing similar series should produce higher compression ratios than when doing so with very different data. This approach appears to be particularly efficient for clustering; it has been applied to fetal heart rate tracings [100]. Following the same underlying ideas, [104] recently proposed a parsing-based similarity distance in order to distinguish healthy patients from hospitalized ones on the basis of various symbolic codings of ECG signals. By comparing the performances of several data classification methods, this distance is shown to be a good compromise between accuracy and computational efforts. Similar approaches have been undertaken earlier in bioinformatics [91] and several compression techniques – such as the Lempel-Ziv complexity [282] – have been successfully applied to compute similarity between biological sequences.

Comparison of distance measures

The choice of an adequate similarity measure highly depends on the nature of the data to analyze as well as application-specific properties that could be required. If the time series are relatively short and visual perception is a meaningful description, shape-based methods seems to be the appropriate choice. If the application is targeting a very specific dataset or any kind of prior knowledge about the data is available, model-based methods may provide a more meaningful abstraction. Feature-based methods seem more appropriate when periodicities is the central subject of interest and causality in the time series is not relevant. Finally, if the time series are long and little knowledge about the structure is available, the compression-based and more generally structure-based approaches have the advantage of being a more generic and parameter-free solution for the evaluation of similarity. Even with these general recommendations and comparisons for the selection of an appropriate distance measure, the accuracy of the similarity chosen still has to be evaluated. Ironically, it seems almost equally complex to find a good accuracy measure to evaluate the different similarities. However (cf. section 5.4.4), a crucial result when indexing is that any distance measure should lower bound the true distance between time series in order to preclude false dismissals [124]. Therefore the tightness of lower bound [209] appears to be the most appropriate option to evaluate the performance of distance measures as it is a completely hardware and implementation independent measure and offers a good prediction concerning the indexing performance. The accuracy of distance measures are usually evaluated within a 1-NN classifier framework. It has been shown by [109] that, despite all proposals regarding different kinds of robustness, the forty year old DTW usually performs better. Table 1 summarizes the properties of every distance measures reviewed in this chapter, based on our formalization of four types of robustness. It also determines whether the distance is a metric and indicates the computational cost and the number of parameters required.

5.4.4 Indexing

An indexing scheme allows to have an efficient organization of data for quick retrieval in large databases. Most of the solutions presented involve a dimensionality reduction in order to index this representation using a spatial access method. Several studies suggest that the various representations differ but slightly in terms of indexing power [209]. However, wider differences arise concerning the quality of results and the speed of

Distance measure	Scale	Warp	Noise	Outlier	Metric	Cost	P
Shape-based							
L_p norms					✓	$O(n)$	0
Dynamic Time Warping (DTW)		✓				$O(n^2)$	1
LB_Keogh (DTW)		✓	✓		✓	$O(n)$	1
Spatial Assembling (SpADe)	✓	✓	✓			$O(n^2)$	4
Optimal Bijection (OSB)		✓	✓	✓		$O(n^2)$	2
DISSIM		✓	✓		✓	$O(n^2)$	0
Edit-based							
Levenshtein				✓	✓	$O(n^2)$	0
Weighted Levenshtein				✓	✓	$O(n^2)$	3
Edit with Real Penalty (ERP)		✓		✓	✓	$O(n^2)$	2
Time Warp Edit Distance (TWED)		✓		✓	✓	$O(n^2)$	2
Longest Common SubSeq (LCSS)		✓	✓	✓		$O(n)$	2
Sequence Weighted Align (Swale)		✓	✓	✓		$O(n)$	3
Edit Distance on Real (EDR)		✓	✓	✓	✓	$O(n^2)$	2
Extended Edit Distance (EED)		✓	✓	✓	✓	$O(n^2)$	1
Constraint Continuous Edit (CCED)		✓	✓	✓		$O(n)$	1
Feature-based							
Likelihood			✓	✓	✓	$O(n)$	0
Autocorrelation			✓	✓	✓	$O(n \log n)$	0
Vector quantization		✓	✓	✓	✓	$O(n^2)$	2
Threshold Queries (TQuest)		✓	✓	✓		$O(n^2 \log n)$	1
Random Vectors		✓	✓	✓		$O(n)$	1
Histogram			✓	✓	✓	$O(n)$	0
WARP	✓	✓	✓		✓	$O(n^2)$	0
Structure-based							
<i>Model-based</i>							
Markov Chain (MC)			✓	✓		$O(n)$	0
Hidden Markov Models (HMM)	✓	✓	✓	✓		$O(n^2)$	1
Auto-Regressive (ARMA)			✓	✓		$O(n^2)$	2
Kullback-Leibler			✓	✓	✓	$O(n)$	0
<i>Compression-based</i>							
Compression Dissimilarity (CDM)		✓	✓	✓		$O(n)$	0
Parsing-based		✓	✓	✓		$O(n)$	0

Table 1: Comparison of the distance measures surveyed in this chapter with the four properties of robustness. Each distance measure is thus distinguished as *scale* (amplitude), *warp* (time), *noise* or *outliers* robust. The next column shows whether the proposed distance is a metric. The cost is given as a simplified factor of computational complexity. The last column gives the minimum number of parameters setting required by the distance measure.

querying. There are two main issues when designing an indexing scheme: *completeness* (no false dismissals) and *soundness* (no false alarms). In an early paper, [124] list the properties required for indexing schemes:

1. It should be much faster than sequential scanning.
2. The method should require little space overhead.
3. The method should be able to handle queries of various lengths.
4. The method should allow insertions and deletions without rebuilding the index.
5. It should be correct, i.e. there should be no false dismissals.

As noted by [213] there are two additional desirable properties:

1. It should be possible to build the index within “reasonable time”.
2. The index should be able to handle different distance measures.

A time series X can be considered as a point in an n -dimensional space. This immediately suggests that time series could be indexed by Spatial Access Methods (SAMs). These allow to partition space into regions along a hierarchical structure for efficient retrieval. B-trees [30] on which most hierarchical indexing structures are based, were originally developed for one-dimensional data. They use prefix separators, thus no overlap for unique data objects is guaranteed. Multidimensional indexing structures – such as the R-tree [31] – use data organized in minimum bounding rectangles (MBR). However, when summarizing data in minimum bounding regions, the sequential nature of time series cannot be captured. Their main shortcoming is that wide MBR produce large overlap with a majority of empty space. Queries therefore intersect with many of these MBRs.

Typical time series contain over thousand datapoints and most SAM approaches are known to degrade quickly at dimensionality greater than 8-12 [81]. The degeneration with high dimensions caused by overlapping can result in having to access almost the entire dataset by random I/O. Therefore, any benefit gained when indexing is lost. As R-trees and their variants are victims of the phenomenon known as the ‘*dimensionality curse*’ [51], a solution for their usage is to first perform dimensionality reduction. The X-tree (extended node tree), for example, uses a different split strategy to reduce overlap [35]. The A-tree (approximation tree) uses VA-file-style (vector approximation file) quantization of the data space to store both MBR and VBR (virtual bounding rectangle) lower and upper bounds [333]. The TV-tree (telescopic vector tree) is an extension of the R-tree. It uses minimum bounding regions (spheres, rectangles or diamonds, depending on the type of L_p norm used) restricted to a subset of active dimensions. However, not all methods rely on SAM to provide efficient indexing. [292] proposed the use of suffix trees [163] to index time series. The idea is that distance computation relies on comparing prefixes first, so it is possible to store every series with identical prefixes in the same nodes. The subtrees will therefore only contain the suffixes of the series. However, this approach seems hardly scalable for longer time series or more subtle notions of similarity. In [124] the authors introduced the GEneric Multimedia INdexIng method (GEMINI) which can apply any dimensionality reduction method to produce efficient indexing. [418] studied the problem of multi-modal similarity search in which users can choose between multiple similarity models depending on their needs. They introduced an indexing scheme for time series where the distance

function can be any \mathcal{L}_p norm. Only one index structure is needed for all \mathcal{L}_p norms. To analyze the efficiency of indexing schemes, [174] considered the general problem of database indexing workloads (combinations of data sets and sets of potential queries). They defined a framework to measure the efficiency of an indexing scheme based on two characterizations: *storage redundancy* (how many times each item in the data set is stored) and *access overhead* (how many unnecessary blocks are retrieved for a query). For indexing purposes, envelope-style upper and lower bounds for DTW have been proposed [211]; the indexing procedure of short time series is efficient but similarity search typically entails more page reads. This framework has been extended [391] in order to index multidimensional time series with DTW as well as LCSS. [13] proposed the TS-tree, an indexing method offering efficient similarity search on time series. It avoids overlap and provides compact meta data information on the subtrees, thus reducing the search space. In [224], the use of an Incremental DFT Computation index (IDC-Index) has been proposed to handle streams based on a deferred update policy and an incremental computation of the DFT at different update speeds. However, the maintenance of the R*-tree for the whole streaming series might cause a constantly growing overhead and the latter could result in performance loss. It is also possible to use indexing methods to speed up DTW calculation; however, it induces a tradeoff between efficiency and I/O cost. However, [351] recently showed that for datasets that are large enough, the benefits of using DTW instead of Euclidean distance is almost null, as the larger the dataset, the higher the probability to find an exact match for any time series. They proposed an extension of the SAX representation – called indexable SAX (iSAX) – allowing to index time series with zero overlap at leaf nodes.

5.5 RESEARCH TRENDS AND ISSUES

Time series data mining has been an ever growing and stimulating field of study that has continuously raised challenges and research issues over the past decade. We discuss in the following open research issues and trends in time series data mining for the next decade.

STREAM ANALYSIS The last years of research in hardware and network research has witnessed an explosion of streaming technologies with the continuous advances of bandwidth capabilities. Streams are seen as continuously generated measurements which have to be processed in massive and fluctuating data rates. Analyzing and mining such data flows are computationally extreme tasks. Several papers review research issues for data streams mining [139] or management [145]. Algorithms designed for static datasets have usually not been sufficiently optimized to be capable of handling such continuous volumes of data. Many models have already been extended to control data streams, such as clustering [110], classification [182], segmentation [214] or anomaly detection [95]. Novel techniques will be required and they should be designed specifically to cope with the ever flowing data streams.

CONVERGENCE AND HYBRID APPROACHES A lot of new tasks can be derived through a relatively easy combination of the already existing tasks. For instance, [243] proposed three approaches, polynomial, DFT and probabilistic, to predict the unknown values that have not fed into the system and answer queries based on forecast data. This approach is a combination of prediction (cf. section 5.3.5) and query by content (cf. section 5.3.1) over data streams. This work shows that future research has to rely

on the convergence of several tasks. This could potentially lead to powerful hybrid approaches.

EMBEDDED SYSTEMS AND RESOURCE-CONSTRAINED ENVIRONMENTS With the advances in hardware miniaturization, new requirements are imposed on analysis techniques and algorithms. Two main types of constraints should absolutely be met when hardware is inherently limited. First, embedded systems have a very limited memory space and cannot have permanent access to it. However, most methods use disk-resident data to analyze any incoming informations. Furthermore, sensor networks (which are frequently used in embedded systems) usually generate huge amounts of streaming data. So there is a vital need to design space efficient techniques, in terms of memory consumption as well as number of accesses. An interesting solution has been recently proposed in [416]. The algorithm is termed *autocannibalistic*, meaning that it is able to dynamically delete parts of itself to make room for new data. Second, as these resource-constrained environments are often required to be autonomous, minimizing energy consumption is another vital requirement. [44] has shown that sending measurements to a central site in order to process huge amounts of data is energy inefficient and lack scalability.

DATA MINING THEORY AND FORMALIZATION A formalization of data mining would drastically enhance potential reasoning on design and development of algorithms through the use of a solid mathematical foundation. [123] examined the possibility of a more general theory of data mining that could be as useful as relational algebra is for database theory. They studied the link between data mining and Kolmogorov complexity by showing their close relatedness. They conclude from the undecidability of the latter that data mining will never be automated, and therefore stating that “*data mining will always be an art*”. However, a mathematical formalization could lead to global improvements of both reasoning and the evaluation of future research in this topic.

PARAMETER-FREE DATA MINING One of the major problems affecting time series systems is the large numbers of parameters induced by the method. The user is usually forced to “fine-tune” the settings in order to obtain best performances. However, this tuning highly depends on the dataset and parameters are not likely to be explicit. Thus, parameter-free systems is one of the key issue that has to be addressed. [216] proposed a first step in this direction by introducing a compression-based algorithm which does not require any parameter. As underlined by [123], this approach could lead to elegant solutions free from the parameter setting problem.

USER INTERACTION Time series data mining is starting to be highly dedicated to application specific systems. The ultimate goal of such methods is to mine for higher-order knowledge and propose a set of solutions to the user. It could therefore seem natural to include an user interaction scheme to allow for dynamic exploration and refinement of the solutions. An early proposal by [210] allows for relevance feedback in order to improve the querying process. From the best results of a query, the user is able to assign positive or negative influences to the series. A new query is then created by merging the series with respect to the user factors on which the system iterates. Few systems have tried to follow the same direction. However, an interactive mining environment allowing dynamic user exploration could increase the accessibility and usability of such systems.

EXHAUSTIVE BENCHMARKING A wide range of systems and algorithms has been proposed over the past few years. Individual proposals are usually submitted together with specific datasets and evaluation methods that prove the superiority of the new algorithm. As noted by [208], selecting those datasets may lead to *data bias* and showed that the performance of time series systems is highly data-dependent. The superiority of an algorithm should be tested with a whole range of datasets provided by various fields [109]. There is still a need for a common and exhaustive benchmarking system to perform objective testing. Another highly challenging task is to develop a procedure for real-time accuracy evaluation procedure. This could provide a measure of the accuracy achieved, thus allowing to interact with the system in real-time to improve its performance.

ADAPTIVE MINING ALGORITHM DYNAMICS Users are not always interested in the results of a simple mining task and prefer to focus on evolution of these results in time. This actually represents the *dynamics* of a time series data mining system. This kind of study is of particular relevance in the context of data streams. [111] studied what are the distinctive features of analyzing streams are, rather than other kinds of data. They argued that one of the core issues is to mine *changes* in data streams. As they are of constantly evolving nature, a key aspect of the analysis of such data is to establish how an algorithm is able to adapt dynamically to such continuous changes. Furthermore, this could lead to ranking changes on the basis of relevance measures and contribute to the elaboration of methods to summarize and represent changes in the system. By finding a way to measure an approximate accuracy in real-time, it should be possible to imagine more “morphable” algorithms that could adapt dynamically to the nature of the data available on the basis of their own performances.

LINK TO SHAPE ANALYSIS Shape analysis has also been matter for discussion over the past few years. There is an astonishing resemblance between the tasks that have been examined; such as query by content [43], classification [205], clustering [246], segmentation [343] and even motif discovery [406]. As a matter of fact, there is a deeper connection between these two fields as recent work shows the numerous inherent link existing between these. [26] studied the problem of classifying ordered sequences of digital images. When focusing on a given pixel, it is possible to extract the time series representing the evolution of the information it contains. As this series is morphologically related to the series of the neighboring pixels, it is possible to perform a classification and segmentation based on this information. As presented above, [417] proposed to extract a time series from the contour of an image. They introduced the time series shapelets that represents the most informative part of an image and allows to easily discriminate between image classes. We can see from these works that both fields could benefit from each other. Even if only modest progress has been made in that direction, a convergence of both approaches could potentially lead to powerful systems.

MULTIOBJECTIVE OPTIMIZATION

Multiobjective approaches were designed to handle problems where several objectives are required to be optimized simultaneously. In order to achieve such a joint optimization, an alternative notion of optimality needs to be adopted. This concept was originally introduced by Francis Edgeworth in [117] and later generalized by the economist Vilfredo Pareto [290] and is called the *Pareto optimum*. As we will detail later, a Pareto solution is optimal in each direction of optimization. Therefore, multiobjective methods can be used to analyze and select between several potentially feasible options. Furthermore, the core strength of the multiobjective approach is that it allows a high degree of flexibility. Indeed, the optimal solution can perform extremely badly on a dimension as long as they perform extremely well on another. The distribution of such solutions is usually referred to as the *Pareto front* [85]. Multiobjective optimization has been applied to many real-world situations like natural resource management [267], medical diagnosis [32] and chemical engineering [45]. We refer the interested reader to recent reviews on multiobjective optimization and analysis techniques Ehrgott [119], Figueira et al. [128].

6.1 DEFINITIONS

Multiobjective minimization problems are defined by a given search space (sometimes called *decision space*) S and a set of functions $F = \{f_1, \dots, f_N\}$ to minimize over S . Formally, a multiobjective problem is defined by

$$\begin{cases} \min & F(x) = \{f_1(x), \dots, f_N(x)\} \\ \text{s.t.} & x \in S \end{cases} \quad (6.1)$$

Definition 19. Let S be a decision space and F the set of functions of a multiobjective problem over S . The *criteria space* is defined as

$$C = \{(f_1(x), \dots, f_N(x)) \mid x \in S\} \quad (6.2)$$

Usually, the ideal solution x_{ideal} , which is the global minimum for all criteria does not exist

$$\nexists x_{\text{ideal}} \in S, \forall n \in \{1, \dots, N\}, f_n(x_{\text{ideal}}) = \min_S f_n \quad (6.3)$$

Therefore, multiobjective problems cannot be solved with a single “perfect” solution, but rather with a *set of efficient solutions* called *Pareto solutions*.

6.1.1 Pareto dominance

As multiobjective optimization is based on finding a set of tradeoffs among potential solutions, we need a relaxed definition of dominance. An efficient solution is, therefore,

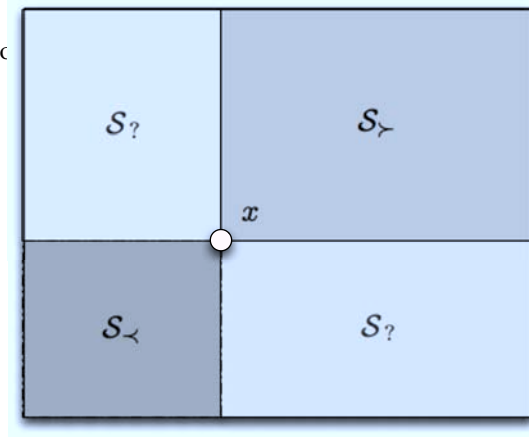


Figure 16: *Pareto dominance* relations for a minimization problem in a bi-criteria space. Any point x of the criteria space divides it into three sub-spaces depending on the dominance relation. $S_{<}$ contains the elements that dominates x ($\forall y \in S_{<}, y \prec x$). Elements of $S_{>}$ are dominated by x ($\forall y \in S_{>}, x \prec y$). Finally, the elements of $S_{?}$ simply cannot be compared to x as they are not dominated nor dominate x ($\forall y \in S_{?}, x \not\prec y \wedge y \not\prec x$).

a solution that is not dominated in *every* objective. In other words, it is impossible to find another solution that jointly improves the complete set of criteria of an efficient solution.

Definition 20. Let x and y be two points of a search space S . We say that a solution y is *Pareto dominated* by a solution x (noted $x \preceq y$) if and only if it is dominated in every dimension. More formally

$$\forall n \in \{1, \dots, N\}, f_n(x) \leq f_n(y) \quad (6.4)$$

In that case y is said to be a *weakly dominated* by x

Definition 21. We say that x *strictly dominates* y (noted $x \prec y$) iff

$$\begin{cases} \forall n \in \{1, \dots, N\}, & f_n(x) \leq f_n(y) \\ \exists n_0 \in \{1, \dots, N\}, & f_{n_0}(x) < f_{n_0}(y) \end{cases} \quad (6.5)$$

We say that a solution y is *strongly dominated* by a solution x (noted $x \prec\prec y$) if and only if it is strictly dominated in every dimension. More formally

$$\forall n \in \{1, \dots, N\}, f_n(x) < f_n(y) \quad (6.6)$$

The dominance relation \prec induces only a partial order on the criteria space, as shown in Figure 16. For any element x , the criteria space is divided into three regions depending on the dominance relation between x and the corresponding subspaces. We call these three subspaces $S_{<}$, $S_{>}$ and $S_{?}$. $S_{<}$ contains the elements that dominate x ($\forall y \in S_{<}, y \prec x$). $S_{>}$ is the subspace whose elements are dominated by x ($\forall y \in S_{>}, x \prec y$). Finally, the elements of $S_{?}$ just cannot be compared to x as they are not dominated nor dominate x ($\forall y \in S_{?}, x \not\prec y \wedge y \not\prec x$).

Definition 22. The set of non-dominated (or efficient) elements of S is called the *Pareto front*.

Solving a multiobjective problem can thus be summarized as the discovery of the Pareto front. Figure 17 depicts a search space in the bi-objective case. We can clearly

see in this figure the Pareto front which emerges from all the non-dominated solutions. From a strict point of view, none of the Pareto solutions can be preferred to others. As we can see, this approach provides a high degree of flexibility, as all directions of optimization are allowed.

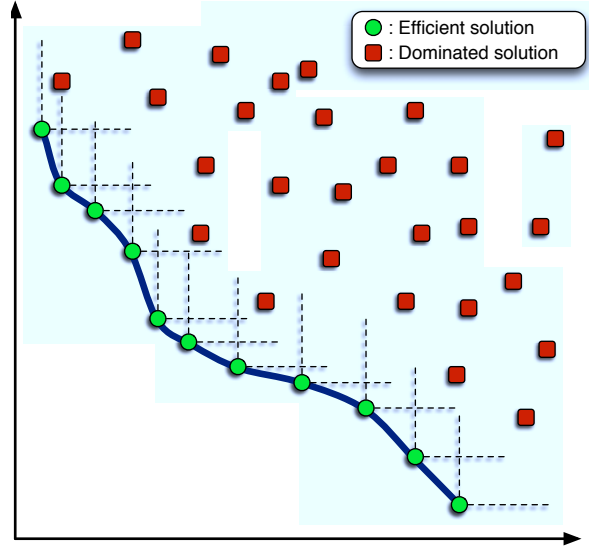


Figure 17: Efficient solutions and dominated solutions for a bi-criteria minimization problem. We can clearly see the *Pareto front* of non-dominated solutions. The dotted lines define the sub-spaces which are dominated by these solutions.

6.1.2 Chebyshev norms

Now, if we look at the problem from the opposite side, each non-dominated element of the Pareto front is the best solution of a monoobjective problem with respect to a set of weights.

Definition 23. Let Λ be the subset of $[0; 1]^N$ such that

$$\forall \lambda = (\lambda_1, \dots, \lambda_N) \in \Lambda, \quad \sum_i \lambda_i = 1 \quad (6.7)$$

The norm defined over a criteria space C by

$$\forall x \in C, \quad \|x\|_\lambda = \max_i \lambda_i x_i \quad (6.8)$$

is called the *weighted Chebyshev norm* given the set of weights $\lambda \in \Lambda$.

Proposition 24. A solution x belongs to the Pareto front if and only if there exists a set of weights $\lambda \in \Lambda$ such that

$$x = \operatorname{argmin}_{y \in C} \|y\|_\lambda \quad (6.9)$$

Figure 18 illustrates this duality between a multiobjective problem and a set of monoobjective problems. The efficient configurations x_a , x_b and x_c , are the best

solutions of an associated monoobjective minimization problem given respectively by norms \mathcal{N}_a , \mathcal{N}_b and \mathcal{N}_c . For each of them, the associated sets of weights defines their weighted Chebyshev norm and, therefore, different directions of optimization.

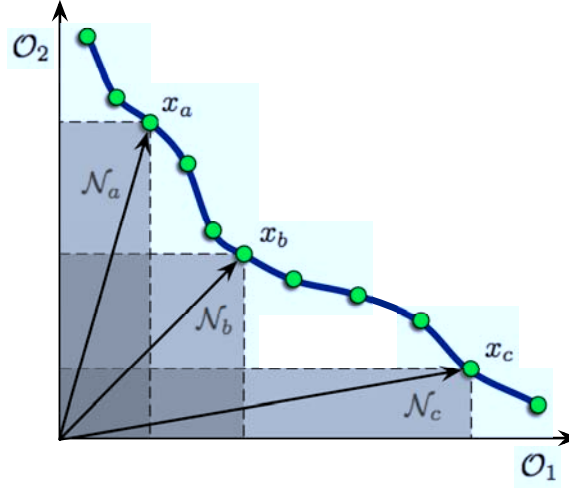


Figure 18: Three efficient solutions x_a, x_b, x_c and their corresponding induced weighted Chebyshev norms $\mathcal{N}_a, \mathcal{N}_b, \mathcal{N}_c$ in a bi-objective problem. Each point of the Pareto front is thus the best solution of a mono-objective problem weighted by its corresponding Chebyshev norm.

6.2 ALGORITHMS CLASSIFICATION

Over the years, various methods and algorithms have been proposed to solve multi-objective optimization problems. Overall, these methods can be divided in two broad categories depending on their approach. First, the *complete methods* allow to find the best optimal solutions (global minimum), but do not provide any limit on their execution time. Therefore, they sometimes can not be feasible depending on the complexity of the problem at hand. Second, the *metaheuristics* allow to obtain approximate solutions in a fixed amount of time. However, the outcome of such algorithms is usually not deterministic and can not provide absolute guarantees on the quality of found solutions. As usually the case with combinatorial problems, we can see that there is a tradeoff between the *efficiency* (execution time) and the *effectiveness* (quality of solutions) of the algorithms.

When the properties of the problem are of limited complexity (small number of objectives, linear objective functions, restricted criteria space), complete methods can provide an optimal solution to multiobjective optimization. In these cases, classical optimization tools such as dynamic programming, the A^* search algorithm or *branch and bound* algorithms are preferred because they guarantee the optimality of solutions as they span the search space entirely. However, for problems that require more than two objectives or with a wide search space, such methods are impossible to use because of the search complexity and combinatorial explosion.

The *metaheuristics* allow to provide a turnaround for these shortcomings by providing approximate solutions in a reasonable amount of time, even if the solutions are usually suboptimal. The effectiveness of such methods can not be proved theoretically but are exhibited through empirical experiments. The proposed metaheuristics for solving

multiobjective optimization problems broadly falls into two distinct categories. First, the *neighborhood* methods try to refine iteratively a single configuration in order to converge towards the Pareto front. Second, the *population* methods use a set of configurations that interact together in order to optimize jointly the overall quality of the set. This approach has been the most studied in the literature, as exemplified by the popularity of *genetic algorithms* (GA) in multiobjective optimization. We redirect the interested readers to reviews on the proposed approaches for multiobjective combinatorial optimization problems Ehrgott and Gandibleux [120].

6.3 APPLICATIONS

Multiobjective optimization has been applied in a wide variety of scientific fields, often towards a system that could help in decision-making processes. Multicriteria methods can thus be used to analyze and select between different potentially feasible water resources development options [3] where it has been applied to many real-world situations [149] based on criteria such as environmental protection, water demand, regional cooperation. Similar lines of research have been followed for conceptual runoff [345] or hydrologic model calibration [415], for fisheries management [260] or environmental decision making [219]. Multicriteria analysis has also been widely used for agricultural resource management [169], either for selecting the optimal multiattribute alternative or for solving multiobjective planning problems. Romero and Rehman [331] underlined the suitability of MCDM techniques to natural resource management. They further presented the theoretical aspects of multicriteria techniques and detailed practical considerations for the management of agricultural systems [326]. A good review of multicriteria decision analysis applied to forest management and other natural resources has been presented in [267].

A lot of research has been devoted to applying multiobjective techniques to assist medical diagnosis. Belacel developed a fuzzy multicriteria classification method called PROAFTN [32] which is a supervised classification scheme. This technique is used to solve the *nominal sorting problematic* [299] where unordered categories are represented by reference objects called *prototypes*. They also showed its applicability to aid diagnosis of stercytic and bladder tumours [33] as well as acute leukemia [34]. Based on a patient's symptoms, Technique Ordered Preference by Similarity to the Ideal Solution (TOPSIS) is one of the widely used methods in medical diagnosis systems [317]. More recently, Zhang et al. [430] developed a linear programming approach for improving the accuracy of medical diagnosis as well as prognosis. Application of multiobjective optimization to chemical engineering has also seen a flourishing literature in the past years [45], with problems related to optimization of chemical processes design [346], polymerization reactions [380], waste treatment [370] and air pollution control [323]. Multiobjective optimization has even been applied to beer dialysis [423] in order to preserve a beer's unique taste even with lowered alcohol content.

Part III

MULTIOBJECTIVE TIME SERIES (MOTS) MATCHING

7

MOTS FRAMEWORK

As we discussed in Section 2.4, studying the properties of musical atoms and auditory perception has taught us that we process audio features in both a temporal and multidimensional manner. As we also discussed in Section 3.3, we are able to perceive uncorrelated features simultaneously and perform flexible similarity assessments. Hence, it seems that a single measure would be unlikely to convey such perceptual similarity (as already pointed out by several authors [115, 114, 386]). Therefore, we believe that both this multidimensional nature and the ability to perceive complex temporal structures shall be taken into account in similarity matching. However, no current retrieval system seems to address these limitations (even outside the realm of audio matching). For example, audio retrieval techniques usually just borrow from other fields such as pattern recognition in order to obtain a single audio similarity measure. While they clearly address many engineering problems, they do not expressly address the multidimensional issues involved in the similarity of timbre. Hence, we wish to incorporate time series information and at the same time seek a multidimensional assessment of similarity. Motivated by these observations, we introduce the generic *MultiObjective Time Series* (MOTS) matching problem. This problem can be applied to any problem in which various time series should be matched jointly, without favoring any dimension in the process. The goal of MOTS matching is therefore to provide a flexible comparison of multiple time series.

Equipped with the basic notions of time series matching and multiobjective optimization, we begin by introducing the general MOTS matching problem (Section 7.1). We indicate the core differences between this novel problem and multivariate matching (Section 7.2) and briefly discuss its complexity. We then introduce two algorithms to solve this problem (Section 7.3) and show their efficiency on massive sets of data (Section 7.4). We compare the efficiency of these algorithms between real and synthetic sets of data (Section 7.4.2). Finally, we discuss the application of the MOTS framework to audio retrieval settings. We show that, based on this framework, we can easily construct two innovative audio querying paradigms (Section 7.5), that allow to go beyond the traditional audio query applications.

7.1 PROBLEM DEFINITION

Problem 25. A *Multiobjective Time Series* (MOTS) matching problem is defined as finding the efficient elements of a database that jointly minimize a set of time series distances

$$\begin{cases} \min \mathcal{D}_Q^k(\mathcal{S}) & k \in \{1, \dots, K\} \\ \text{s.t. } \mathcal{S} \in \text{DB} \end{cases} \quad (7.1)$$

with Q the query represented by a set of K time series and \mathcal{S} the elements of the database DB which contains time series corresponding to the same objectives as the query. Finally, $\mathcal{D}_Q^k(\mathcal{S})$ is the similarity between the k^{th} feature represented by time series Q_k and \mathcal{S}_k , i.e. $\mathcal{D}_Q^k(\mathcal{S}) = \mathcal{D}(Q_k, \mathcal{S}_k)$ (cf. Definition 5).

It is necessary to note here that this problem is not a problem of *optimization* (as are usual multi-objective methods). Indeed, the elements in the database are fixed and, therefore, there is no *feature space*. However, we can already see now that part of the computational complexity of this problem arises from objective functions $\mathcal{D}_Q^k(\mathcal{S})$ which represent time series distances. As we discussed in Section 5.4.3, time series similarity is a remarkably subtle concept that can entail a high computational complexity. Furthermore, because of the multiobjective nature of this problem, it is impossible to gain straightforward efficiency from traditional time series indexing methods. Indeed, these techniques provide most of their pruning power by avoiding computation of irrelevant parts of the search space. As we will discuss further in Section 7.2, it is noteworthy here to understand the fundamental differences between multivariate time series matching (extensively studied in literature) and our multiobjective problem. Multivariate problems usually imply that the series are somehow statistically linked and attempt to find a single similarity measure to compare a *set* of time series. This allows to circumvent the problem of pruning power raised by the notion of Pareto dominance, which is the second aspect of computational complexity in the MOTS problem. Indeed, multiobjective problems allow the optimization of objectives that can be conflictive with each other. Finding the most similar item \mathcal{S}^* to a MOTS query requires to jointly minimizing the distances between two sets of time series.

$$\mathcal{S}^* = \underset{\mathcal{S}}{\operatorname{argmin}} \left\{ \left(\mathcal{D}_Q^k(\mathcal{S}) \right), k = 1, \dots, K \right\} \quad (7.2)$$

As the *ideal point* \mathcal{S}^* which simultaneously optimizes all criteria usually does not exist, solving this problem turns out to finding the set of *tradeoff solutions* that offer different compromises among objectives. A solution \mathcal{S} is optimal if there is no other solution in the search space that achieves similarities higher than \mathcal{S} on *every* criterion $\mathcal{D}_Q^k(\mathcal{S})$. Therefore, we can not just rule out a portion of the database if it performs poorly on one objective, as it could contain the best element in another objective. This implies that if we want to know which elements belong to the exact Pareto front, we should evaluate the complete set of distances for every objective, thus degenerating to *brute force* analysis. Figure 19 illustrates these concepts. The query is a set of time series input to the system. The query is at the origin of the criteria space as distances with itself are null in every objective. There is no element in the database that perfectly matches these two time series. Solution \mathcal{A} is the best match for objective \mathcal{O}_1 . As we can see, its first time series is closely similar to that of the query. Solution \mathcal{B} is respectively the best match for objective \mathcal{O}_2 . Finally, element \mathcal{C} is the best solution for the associated mono-objective problem with equal weights. We can see that it is not closely similar to any objective, which exhibits the relevance of our approach. Indeed, the MOTS matching allows joint queries on several dimensions without favoring any of them during the search. Therefore, this approach is an appropriate model when the relative weights of each objective cannot be known in advance, which is particularly relevant for audio perception (cf. Section 5). Depending on the problem, regions of the Pareto front might be preferred to others, according to personal preferences as we will demonstrate in Section 8.4.1. Finally, it is already appealing to note here that the objective functions $\mathcal{D}_Q^k(\mathcal{S})$ can be defined differently depending on the feature being studied. Therefore, the distance function in each dimension can be tailored to fit its corresponding feature. It is even possible to assign multiple objectives for the same time series features, with each dimension corresponding to a different measure of similarity over the same feature.

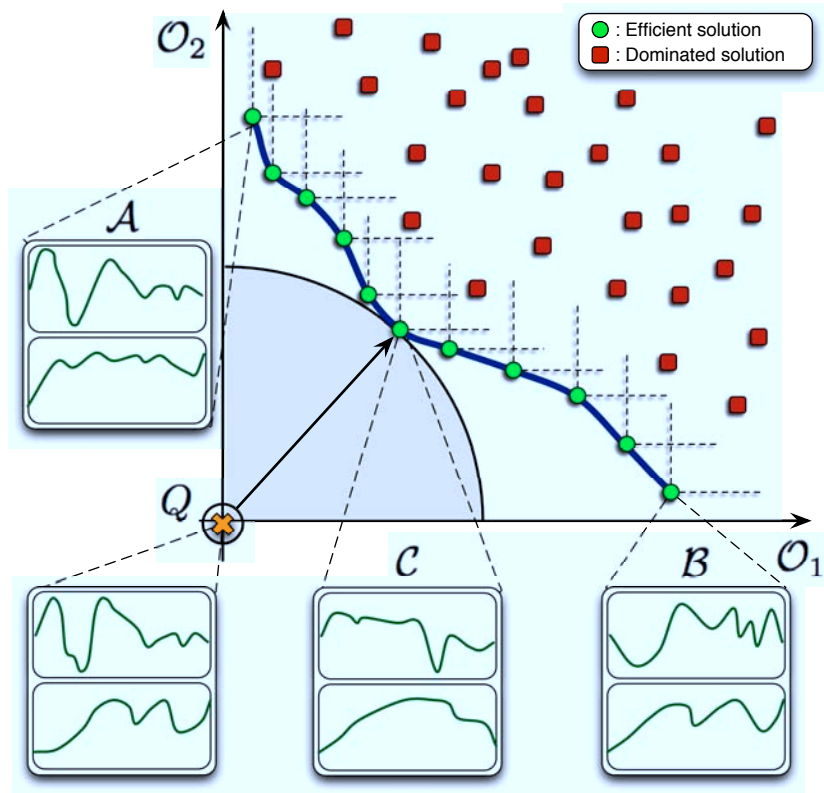


Figure 19: Illustration of the MOTS matching problem in a bi-objective context. The query Q is at the origin of the space and is represented by a set of time series that have to be matched jointly. Solution A is the best match for objective O_1 , as we can see the first time series is closely similar to that of the query. Solution B is respectively the best match for objective O_2 . The element C would be the best solution for the weighted monobjective problem given a set of equal weights. We can see that it is not closely similar to any objective, which motivates the use of multiobjective optimization.

7.2 COMPARISON TO MULTIVARIATE MATCHING

We try to demonstrate the fundamental differences between a multivariate matching algorithm (extensively studied in the literature) and our multiobjective problem. First, as we said earlier, multivariate analysis generally relies on the premise that the set of time series is somehow statistically linked. Sometimes, inference is even possible between the different series of the set. In this manner, multivariate analysis provides mechanisms to obtain trends *across* multiple dimensions and take into account the effect of all variables observed. Oppositely, multiobjective matching allows the optimization of objectives that can be uncorrelated and even conflictive with each other. As we have seen in Section 3.3, we can *perceive* conflictive temporal evolutions, which makes multiobjective approaches even more appealing. Furthermore, multivariate matching usually tries to find a single measure that could encapsulate the similarity on the *whole set* of time series. In this line of thought, a multivariate search can be seen as equivalent to a mono-objective problem spanning several dimensions. Indeed, the multivariate case can be seen as a reduction to a weighted multiobjective search. Hence, weighting and merging the objectives could allow to circumvent the problem of pruning power raised by the notion of Pareto dominance. Finally, the number of dimensions of a multivariate problem is usually fixed whereas the multiobjective approach can work on different subsets of objectives. We illustrate these concepts by trying to find the solution of the similarity problem exposed in Figure 5 with a multivariate nearest-neighbor approach. We can see in Figure 20 that the computation of the Euclidean distance (\mathcal{L}_2) on each feature gives a slight difference in results (because of the small oscillating segment in *loudness*). If we try to find which elements are similar to \mathcal{S}_2^1 , the system will rate element \mathcal{S}_1^1 as being the most similar and then element \mathcal{S}_2^2 as being less similar. Therefore, this approach impose an implicit preference towards the *pitch* of different sounds in similarity matching.

Now we try to solve the same similarity problem with a MOTS approach in Figure 21. This time, if we try to find which elements are more similar to \mathcal{S}_2^1 , the system will isolate the problem depending on its two underlying dimensions. Therefore, element \mathcal{S}_1^1 and \mathcal{S}_2^2 are selected as *efficient* as they are not dominated in any dimension. However, there is no ranking between these two elements and they are treated as being “equally efficient” to the similarity towards \mathcal{S}_2^1 . Only element \mathcal{S}_1^2 is clearly exhibited as being least similar, as it is dominated by the others. Therefore, there is no imposed preference towards any dimension of different sounds in similarity matching. The MOTS matching treats the two dimensions separately and equally.

7.3 ALGORITHMS

Because of the ever-growing size of storage capacities, linear scan of an entire database has become unacceptable. Hence, it would be highly desirable to obtain a search method with sublinear time complexity. We introduce two algorithms that can handle the MOTS matching problem. However, because of the novelty of this approach, no competing method exists to evaluate the efficiency of our algorithms. Therefore, the *multiobjective brute force* algorithm will be our testing baseline. This approach requires to compute every distance in each objective and then extract the Pareto front from the full distance matrix. We describe this reference procedure in Algorithm 7.1.

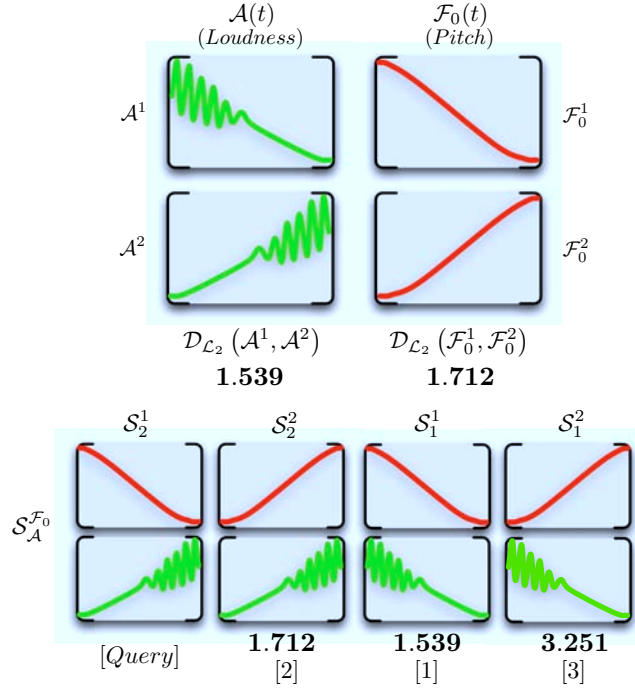


Figure 20: Trying to find the solution to the similarity problem exposed in Figure 5 with a multivariate nearest-neighbor approach. The system will order element \mathcal{S}_1^1 as being the most similar to \mathcal{S}_2^1 , as the distance is slightly different between the two features. Therefore, there is an implicit preference towards the *pitch* of different sounds in similarity matching.

7.3.1 Multiobjective early abandon

As we discussed in Section 7.1, the complexity of the MOTS problem essentially lies in the repeated computations of time series distances. A natural idea would be, therefore, to find a way to restrict the amount of distance computations. Instead of computing the distances for every series and each objective, we would like to drop calculations as soon as we are confident that the corresponding element is dominated. This technique is known as *early abandon*. However, we have to make fundamental modifications in order to account for the multiobjective nature of our problem. Indeed, early abandon in a mono-objective setting is based on comparing the current similarity against the best distance known so far. However, in a multidimensional context where we seek a set of efficient solutions, we cannot simply compare the current distances to a single reference. Therefore, a first turnaround would be to maintain a current working Pareto front with which to compare the successive elements. However, this approach would require to perform several verifications of Pareto dominance at each step of a distance computation. Unfortunately, this verification is a computationally intensive operation. Therefore, it would be preferable to obtain an *approximate distance* for every elements beforehand. That way, we could perform the Pareto verification on these approximations and only compute the complete distances of potentially

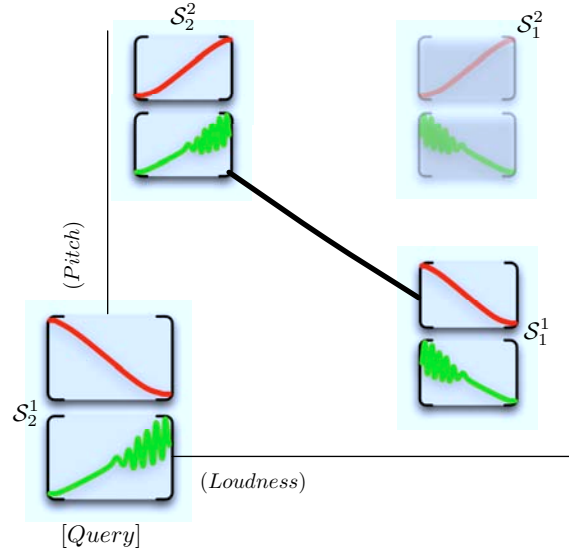


Figure 21: Trying to find the solution of the similarity problem exposed in Figure 5 with a MOTS approach. If we seek to find which elements are more similar to S_2^1 , the system will divide the problem in its two inherent dimensions. Therefore, element S_1^1 and S_2^2 are selected as being the most similar as they are not dominated. Only element S_1^2 is exhibited as being least similar. Therefore, there is no implicit preference towards any dimension in similarity matching.

Algorithm 7.1 Brute force multiobjective time series matching algorithm

```

multiobjectiveBruteForce(Q, db)
  for  $i \in [1 \dots \text{size}(db)]$ 
    for  $k \in [1 \dots N_{obj}]$ 
      compute  $\mathcal{D}_Q^k(S_i)$ 
    end
  end
   $\mathcal{P} \leftarrow \text{extractParetoFront}(\mathcal{D}_Q(S_j))$ 
end

```

efficient solutions. The approximate distance should be *lower-bounding*, ie. it should underestimate the true distance.

$$\mathcal{D}_{approx}^k(S_i) \leq \mathcal{D}_{true}^k(S_i) \quad \forall k \in [1, \dots, N_{obj}] \quad (7.3)$$

With this property, we can prune elements as we are sure that they can only perform *worse* than their current position. In simpler words, if a set of lower-bounding distances is dominated, then we are sure that the corresponding set of true distances is dominated. Therefore, we need to obtain *simplified representations* for the collection of time series that can provide a more efficient distance computation. If these representations are coarse enough, they can account for several time series at the same time. In order to obtain such properties, we can use the SAX representation Lin et al. [251] that performs a temporal and amplitude quantification of the series. In this model, the series are first

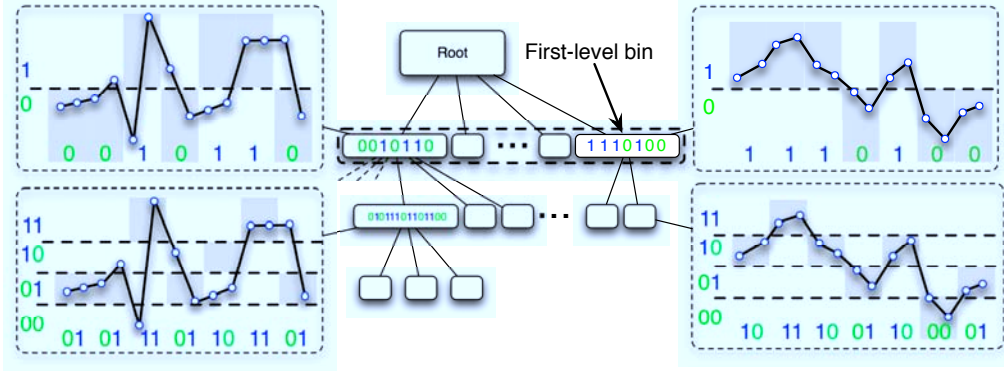


Figure 22: Construction of the quantified bins for time series and computation of the 1st-level distance for a query time series.

divided into a set of equal-sized temporal steps. Then, the average of the time points contained in each step i is computed and matched to an alphabet of reduced size.

$$\bar{\mathcal{T}}_i = \alpha \left(\frac{w}{n} \sum_{j=\frac{n}{w}(i-1)+1}^{\frac{n}{w}i} \mathcal{T}_j \right) \quad (7.4)$$

with n the length of the original series, w the number of resulting temporal steps ($w \ll n$) and $\alpha(x)$ a function that matches $x \in \mathbb{R}$ to a discrete alphabet (amplitude quantification). Based on this representation, the iSAX index [351] provides an efficient tree-like index for time series. The idea behind this index is that each level of the tree provides a finer representation of the series, by increasing the size of the amplitude alphabet. Figure 22 illustrates this construction. The series are divided into 8 equal-sized temporal steps. At the first level, the series are quantified by using an alphabet of two elements $\{0, 1\}$. Then, at the subsequent levels, the series are refined by using a larger alphabet $\{00, 01, 10, 11\}$. Obviously, each node in the tree accounts for a whole set of time series from the database. Hence, if we take the first-level of this representation, we obtain a set of *prototypical bins* of reduced cardinality for the complete database.

Then, the lower-bounding distance between a query Q and a bin representation $\bar{\mathcal{B}}_x$ can be obtained by first transforming the query into the same representation \bar{Q} and then computing

$$\mathcal{D}_{\text{approx}}(\bar{Q}, \bar{\mathcal{B}}_x) = \sqrt{\frac{n}{w}} \sqrt{\sum_{i=1}^w (\mathcal{D}(q_i, b_i))^2} \quad (7.5)$$

Hence, with this construction, we can obtain the lower bounding position of every element in the database, as illustrated in Figure 23. At first glance, it would seem tempting to use these approximate distances to perform a direct assessment of Pareto efficiency. However, it is important to understand that these distances are just lower-bound *approximations*. Therefore, the true dominance relations are still uncertain. This is exhibited in Figure 23 with an outlined relationship. It turns out that the final distance of the potentially dominating element is much higher. However, when its true distances are computed, we are sure that it dominates some of the approximate positions.

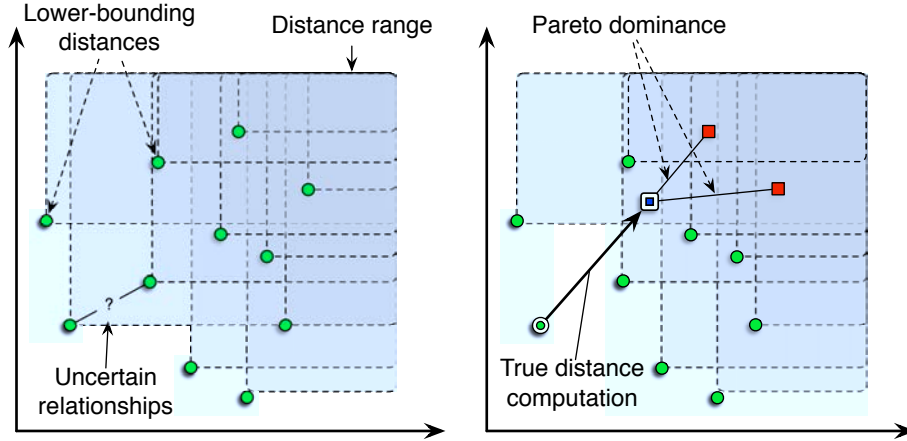


Figure 23: The approximate lower bounding distances in the criteria space and a set of relationships that can or can not be computed

The final implementation is presented in Algorithm 7.2. We start by transforming the query into the quantified representation. Then, we compute the first level distances for all bins. We store these distances for corresponding elements in the distance matrix αDist . We then create an empty Pareto front \mathcal{P} and iterate over the elements of the database. When evaluating an element, as soon as it is dominated by the current Pareto front, we abandon computations. If all the distances have been computed, then the current element is potentially efficient. Therefore, we add this element to the current Pareto front. We compute the new front by removing the eventual dominated points (as the newly added item might dominate existing solutions in the current front). When all the elements of the database have been verified, \mathcal{P} contains the final Pareto front.

7.3.2 Hyperplane search

One of the main problem of the previous algorithm is that it still requires frequent verifications of the Pareto optimality. Hence, it would be wiser to find a less expensive theoretical limit to drop computations of the distance measures. Therefore, our main idea is to construct an approximate Pareto hyperplane \mathcal{P} to act as our theoretical limit. We can obtain this hyperplane by using 1-NN queries from efficient time series indexing for each objective (such as the iSAX index presented in the previous section). These queries will give us boundary elements of the final Pareto front. This is straightforward from the fact that these elements cannot be dominated as they have the smallest distance in one of the objectives.

$$\forall \mathcal{S}_i, \exists k \mid \forall \mathcal{S}_j, \mathcal{D}_Q^k(\mathcal{S}_i) < \mathcal{D}_Q^k(\mathcal{S}_j) \Rightarrow \mathcal{S}_i \in \mathcal{P} \quad (7.6)$$

Hence, we can prune elements whose approximate distances are dominated by this hyperplane. This can be computed straightforwardly if we obtain the hyperplane normal. We show how to compute this normal efficiently by avoiding an expensive least-squares minimization. The normal of a hyperplane can be defined in the following manner

Algorithm 7.2 MOTS matching algorithm with early abandon

multiobjectiveEarlyAbandon(Q, db, idx)

```

// Quantify the query
 $\bar{Q}^{k \in [1 \dots N_{obj}]} = \left\{ \frac{w}{n} \sum_{j=\frac{n}{w}(i-1)+1}^{\frac{n}{w}i} Q_j^k, i \in [1 \dots w] \right\}$ 
// Compute query-to-bin distances
for  $k \in [1 \dots N_{obj}]$ 
  for  $b \in [1 \dots N_{bins}^k]$ 
     $aDist_{i \in \mathcal{B}_b^k}^k = \mathcal{D}_{approx}(\bar{Q}^k, \bar{\mathcal{B}}_b^k)$ 
  end
end
 $\mathcal{P} = \emptyset$ 
// Perform multiobjective abandon
for  $i \in [1 \dots size(db)]$ 
  for  $k \in [1 \dots N_{obj}]$ 
    if isDominated( $aDist_i, \mathcal{P}$ )
      abandon;
    else
       $aDist_i^k = \mathcal{D}_Q^k(S_i)$ 
    end
    add( $S_i, \mathcal{P}$ );
     $\mathcal{P} = extractParetoFront(\mathcal{P});$ 
  end
end

```

Proposition 26. Given a nonzero vector \mathbf{n} in \mathbb{R}^m and a point $\mathbf{p} \in \mathbb{R}^m$, the hyperplane perpendicular to \mathbf{n} through \mathbf{p} is the set of all $\mathbf{x} \in \mathbb{R}^m$ such that

$$(\mathbf{x} - \mathbf{p}) \cdot \mathbf{n} = 0 \quad (7.7)$$

Therefore, if we want to find the normal of hyperplane \mathcal{H} , we must find the vector $\mathbf{n}_p \in \mathbb{R}^m$ satisfying

$$\mathcal{P}\mathbf{n}_p = \mathbf{0}_m \quad (7.8)$$

where \mathcal{P} is a $k \times m$ matrix and $\mathbf{0}_m$ is a $m \times 1$ zero vector. $\mathcal{P} = [p_1, \dots, p_k]$ is the set of Pareto points defining the hyperplane \mathcal{H} (in our case p_i will be the 1-NN result for the i^{th} objective). In order to obtain this vector, we must solve

$$\mathbf{n}_p = \underset{\mathbf{v}}{\operatorname{argmin}} \left(\mathbf{v}^T \mathcal{P}^T \mathcal{P} \mathbf{v} \right) \quad (7.9)$$

Alone, this equation yields the trivial solution $\mathbf{n}_p = \mathbf{0}_m$ which we obviously want to avoid. To avoid this case, we can add the constraint $\|\mathbf{n}_p\| = 1$, which can be rewritten as $1 - \mathbf{n}_p^T \mathbf{n}_p = 0$. Therefore, in order to find the best value for \mathbf{n}_p , we can use the Lagrange multipliers and solve

$$\frac{\delta}{\delta \mathbf{n}_p} \left(\mathbf{n}_p^T \mathcal{P}^T \mathcal{P} \mathbf{n}_p + \lambda \left(1 - \mathbf{n}_p^T \mathbf{n}_p \right) \right) = 0 \quad (7.10)$$

After applying the derivation, we obtain the characteristic equation $(\mathcal{P}^T \mathcal{P} - \lambda E) \mathbf{n}_p = 0$. Therefore, we know that \mathbf{n}_p is an eigenvector of $(\mathcal{P}^T \mathcal{P})$ and λ is an eigenvalue.

However, we can not control the orientation of the normal (as any hyperplane possess two oppositely oriented normal vectors). Furthermore, this also requires to compute some eigenvectors with potentially large dimensionality which can be expensive. In order to alleviate both problems at the same time, we have to slightly modify the original constraint. For that purpose, we introduce a direction vector \mathbf{d} that will ensure the orientation of the normal vector. Therefore, as we constrain the normal vector \mathbf{n}_p to have the same orientation as \mathbf{d} . We can write this constraint as $(1 - \mathbf{d}^T \mathbf{n}_p)^2 = 0$. Hence, we must now solve

$$\mathbf{n}_p = \underset{\mathbf{v}}{\operatorname{argmin}} \left(\mathbf{v}^T \mathcal{P}^T \mathcal{P} \mathbf{v} + (1 - \mathbf{d}^T \mathbf{v})^2 \right) \quad (7.11)$$

By using the same reasoning than previously, we can find the extreme value by solving

$$\frac{\delta}{\delta \mathbf{n}_p} \left(\mathbf{n}_p^T \mathcal{P}^T \mathcal{P} \mathbf{n}_p + (1 - \mathbf{d}^T \mathbf{n}_p)^2 \right) = 0 \quad (7.12)$$

Therefore, by taking the same matrix derivatives and simplifying, we obtain the normal by computing

$$\mathbf{n}_p = \left([\mathcal{P}, \mathbf{d}] [\mathcal{P}, \mathbf{d}]^T \right)^{-1} \mathbf{d} \quad (7.13)$$

where $[\mathcal{P}, \mathbf{d}]$ is the matrix obtained by concatenating matrix \mathcal{P} . In our implementation, we use $\mathbf{d} = \max_j (\mathbf{p}_i^j)$, $\mathbf{p}_i \in \mathcal{P}$ to ensure the orientation of the resulting normal. The distance of any point \mathbf{a}_x relative to the approximate Pareto front \mathcal{P} is then defined as

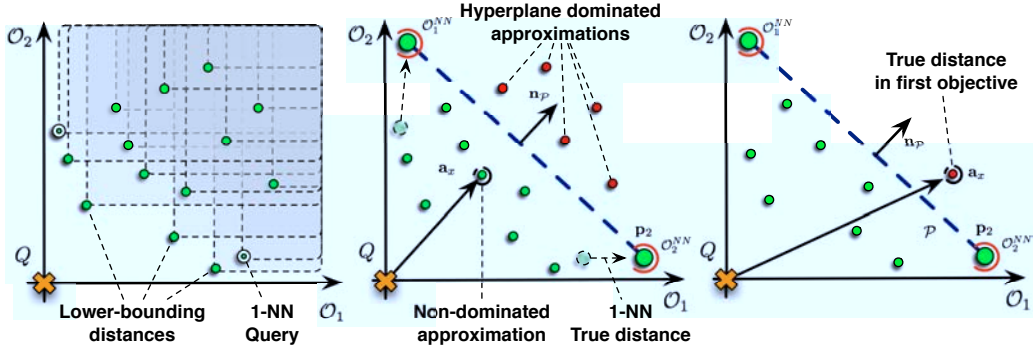
Proposition 27. *Let \mathcal{P} be the hyperplane of all $\mathbf{x} \in \mathbb{R}^k$ with $(\mathbf{x} - \mathbf{p}) \cdot \mathbf{n} = 0$ such that $\mathbf{n} \neq \mathbf{0}$. Then the distance of any point $\mathbf{a}_x \in \mathbb{R}^k$ from the hyperplane \mathcal{P} is given by*

$$\operatorname{dist}(\mathbf{a}_x, \mathcal{P}) = \frac{(\mathbf{a}_x - \mathbf{p}_i) \cdot \mathbf{n}}{\|\mathbf{n}\|} \quad (7.14)$$

with $\|\mathbf{n}\|$ the norm of the normal \mathbf{n} and $\mathbf{p}_i \in \mathcal{P}$ is one of the Pareto points.

The final algorithm is illustrated geometrically in Figure 24. Even if we use the iSAX index, the implementation presented here can be used with any representation, distance and indexing techniques available (cf. Chapter 5). We simply assume that a time series index is constructed for each objective in order to perform efficient 1-NN queries and, therefore, avoid linear scan. We also consider that the index provides a lower bounding distance measure on indexing nodes (as explained in the previous section).

This implementation is presented in Algorithm 7.3. Given a query Q , a database db and a set of index TS-Indexes for each objective (constructed prior to the search), we start by transforming the query and computing the first-level distances as previously. This set aDist is then used to perform the 1-NN exact queries on each objective. These queries give us the initial Pareto front \mathcal{P} that form the approximate Pareto hyperplane. The 1-NN queries also compute a small portion of exact distances for each objective that we recover in list aDist . That way, after 1-NN queries we already have an approximate

Figure 24: Geometric interpretation of the *multiobjective hyperplane search* algorithm

lower bounding position for each element. We then obtain the normal of the hyperplane defined by the list of Pareto points. Then, we evaluate each element of the database and stop distance computation as soon as they are dominated by the hyperplane. If we compute the complete distances in every objective for an element, we add it to the list of potential Pareto points. Finally, we filter this list by extracting the final Pareto front \mathcal{P} at the end of the algorithm.

7.4 EFFICIENCY ON MASSIVE DATABASES

We present the results of our algorithms regarding computation efficiency. Unfortunately, because of the novelty of this problem, there exist no competing method to compare. Hence, we are evaluating our methods against the *brute force multiobjective* algorithm on synthetic and real datasets. The artificial dataset is composed of random walk time series generated with a constant size of 512 time points. An independent set is synthesized for each hypothetical objective. The second (real) dataset is a combination of *Studio On Line* [24], *Real World Computing* [152] and *Vienna Symphonic Library* instrumental databases. These datasets include single notes of different playing modes from 23 orchestral instruments, which amounts to a total of 213.814 sound files. These files are WAVE and AIFF format, quantified to 16-bit at a sampling rate of 44.1 kHz. Subsets of the collections are randomly selected for increasing database sizes. Objectives are also randomly selected from a set of audio descriptors (cf. Table 2). This selection procedure is ten-folded. For each set of parameters (database size and set of objectives), one hundred queries are processed in order to avoid statistical anomalies. Queries are random walk time series with a constant size of 512 points. Computations were performed on a Macbook 2.4 GHz Dual Core running under Mac OS X 10.6.6 with 2 GO of DDR3 RAM.

7.4.1 Comparing algorithms

We present the results of different algorithms in terms of *querying wall time* for synthetic datasets in Figure 25. The left figure shows the *median* (dotted line), *average* and *variance* (solid line) in querying time for increasing database sizes. As we can see, the early abandon algorithm can already provide up to two times of speedup over the brute force approach, with a very low variance in the querying time. However, this factor of speedup appears to be linear to the cardinality of the dataset. The *hyperplane* algorithm

Algorithm 7.3 MOTS matching algorithm by approximate hyperplane search.

```

multiobjectiveHyperplaneSearch(Q, db)
  // Quantify the query
   $\bar{Q}^{k \in [1 \dots N_{\text{obj}}]} = \left\{ \frac{w}{n} \sum_{j=\frac{n}{w}(i-1)+1}^{\frac{n}{w}i} Q_j^k, i \in [1 \dots w] \right\}$ 
  // Compute query-to-bin distances
  for k  $\in [1 \dots N_{\text{obj}}]$ 
    for b  $\in [1 \dots N_{\text{bins}}^k]$ 
       $\text{aDist}_{i \in \mathcal{B}_b^k}^k = \mathcal{D}_{\text{approx}}(\bar{Q}^k, \bar{\mathcal{B}}_b^k)$ 
    end
  end
  // Perform efficient 1-NN queries
  [ $\mathcal{P}$  aDist] = 1NN-Queries(Q, aDist, TS-Indexes)
  // Reference direction vector
   $\mathbf{d} = \max_j(p_i^j), p_i \in \mathcal{P}$ 
  // Compute hyperplane normal
   $\mathbf{n}_p = \left( [\mathcal{P}, \mathbf{d}] [\mathcal{P}, \mathbf{d}]^T \right)^{-1} \mathbf{d}$ 
  // Transform into unit-norm vector
   $\mathbf{n}_p = \mathbf{n}_p / \sqrt{\mathbf{n}_p^T \mathbf{n}_p}$ 
  for i  $\in [1 \dots \text{size}(\text{db})]$ 
    for k  $\in [1 \dots N_{\text{obj}}]$ 
      if  $(\text{aDist}_i - p_1) \cdot \mathbf{n}_p < 0$ 
        abandon
      else
         $\text{aDist}_i^k = \mathcal{D}_Q^k(\mathcal{S}_i)$ 
      end
    end
    add( $\mathcal{S}_i$ ,  $\mathcal{P}$ )
  end
  checkParetoFront(pPoints)

```

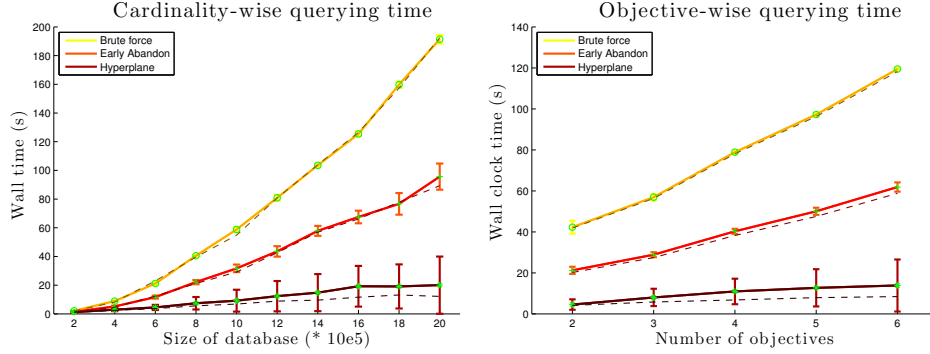


Figure 25: Query wall time (in seconds) for increasing database size (left) and increasing number of objectives (right) on synthetic datasets.

is strongly superior, as it provides up to ten times of speedup over the brute force approach, and its median is usually even faster. The differences between the early abandon and hyperplane approaches can be explained by the higher number of Pareto front evaluations in the first one. However, the variance in querying times of the hyperplane search also increase with the cardinality, which imply that the resulting time might vary more importantly than early abandon depending on the distribution of the dataset. The most enthralling finding concerns the efficiency of our algorithm with increasing number of objectives, presented in Figure 25 (right). As we can see, the early abandon also provide a linear factor of speedup over the brute force approach. However, the hyperplane algorithm exhibits a sub-linear behavior when the number of objectives grows, once again with a significantly lower median. This sub-linear behavior could be explained by the higher probability that a large portion of the search space is ruled out by the approximate hyperplane with a greater number of dimensions.

To analyze this hypothesis, we compare the pruning power induced by each algorithm. The *space pruning ratio* is computed by comparing the proportion of points that are not entirely evaluated (ie. their complete distances are not calculated) to the quantity of points in the dataset. One of the main advantage of this measure is that it is hardware and dataset independent, furthermore it is also independent of the complexity of the distance measure used in the final computation. Therefore, the gain provided by the different techniques can be compared objectively. Figure 26 (left) exhibits the space pruning ratio provided by the early abandon and the hyperplane algorithms (the brute force is omitted as its pruning ratio is obviously null) for a growing amount of time series in the database. As we can see in this figure, the *hyperplane* limit provides a strongly superior pruning ratio as compared to the *early abandon* technique. The variances seem to remain constant (with small variations) for both algorithms as the number of objective grows, with once again a higher variance for the hyperplane algorithm. However, it must be understood that this measure is, in fact, a ratio of the number of elements evaluated. Therefore, an equivalent variance for higher cardinality will imply a higher variance in the number of elements pruned. In both case, the techniques seem to indicate an upper bound in pruning power as the number of time series increase. Figure 26 (right) exhibits the space pruning ratio provided for a growing number of objectives. As we can see, the hyperplane algorithm quickly converge to a constant pruning ratio around 80% (which can explain its sub-linear time complexity), whereas the early abandon algorithm seems to exhibit a continuous drop in pruning power.

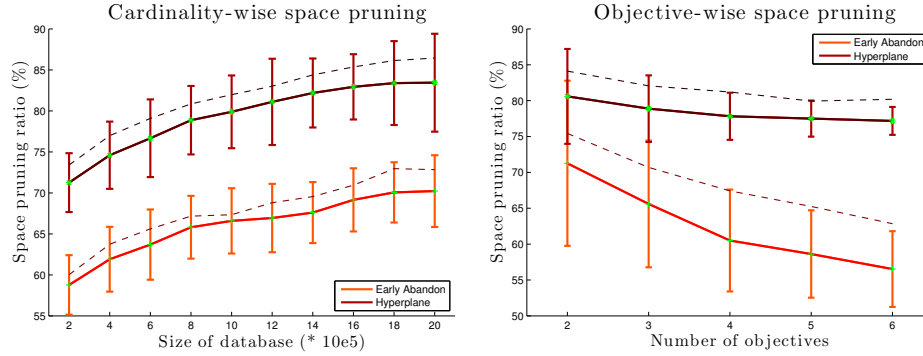


Figure 26: Space pruning ratio for increasing database size (left) and increasing number of objectives on synthetic datasets.

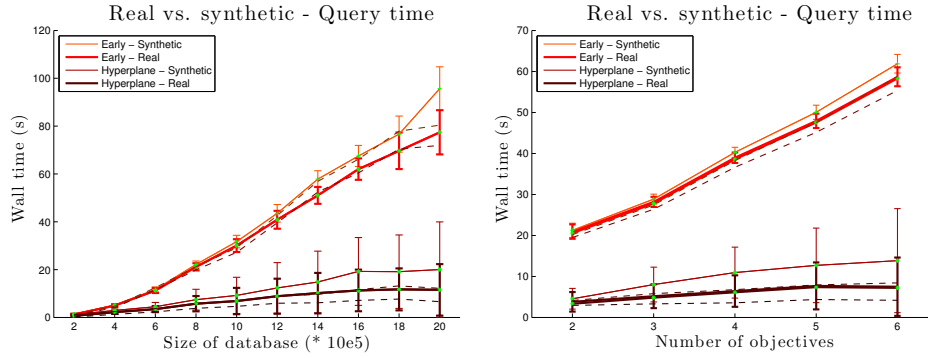


Figure 27: Query wall time (in seconds) for increasing database size (left) and increasing number of objectives (right) compared between synthetic and real datasets.

7.4.2 Comparing datasets

We now compare the performances of different algorithms depending on the nature of underlying data. Even if the comparison on artificial datasets shows the strong superiority of the proposed hyperplane algorithm, a proper evaluation should rely preferably on real datasets. Therefore, we analyze the efficiency of different algorithms on the audio collection presented previously and compare it to the results obtained on synthetic datasets. In order to achieve a meaningful comparison, the time series in the real dataset have all been resampled to a length of 512 time points. We present the results of algorithm wall time speed in Figure 27 for a growing number of series and objectives. We omit the results of brute force for the sake of clarity. Analysis of results reveals that both algorithms performs even better on real sound collections. This could be explained by the distribution of time series in real datasets, which is unlikely to be uniform as it is for random walk datasets. However, it seems that the enhancement is more pronounced for the hyperplane algorithm in both cardinality and objective-wise results. This seems to follow the intuition that the use of an approximate hyperplane benefits from the uneven distribution of data. Therefore, it enhances the overall efficiency of the algorithm.

We offer the same comparison for the space pruning ratio in Figure 28. As we can see, the higher performance for both algorithms is clearly seen in their respective pruning

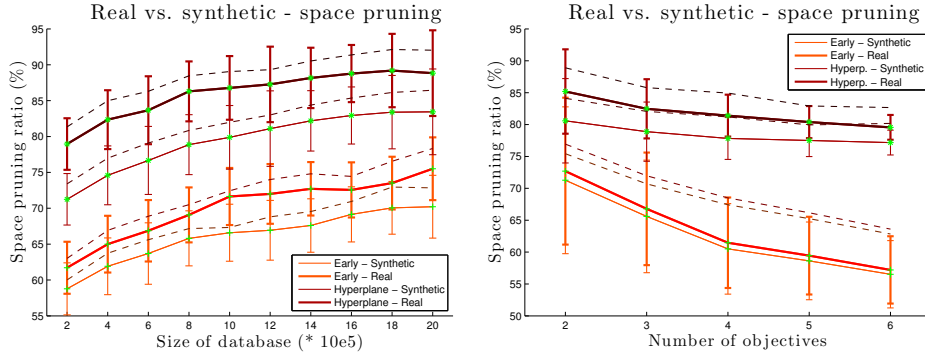


Figure 28: Space pruning ratio for increasing database size (left) and increasing number of objectives (right) on synthetic and real datasets.

power for increasing databases sizes (left). The somehow more chaotic distribution of space pruning ratios can be explained by the nonuniform distributions of data, which cause a wider disparity in each evaluation. For an increasing number of objectives (right), the improvement seems to be a lot less noticeable, with the same overall evolution of pruning power with the number of objectives. In the real dataset, even if the pruning power of the hyperplane method is at first higher, its loss is slightly more substantial with a higher number of objectives.

7.5 INNOVATIVE AUDIO QUERYING

Now equipped with a flexible matching framework based on observation of the auditory perception, it seems logical to apply it to audio settings. We begin by showing the potential interest for our framework in content-based audio retrieval, by outlining the limitations in state-of-art audio matching methods (Section 7.5.1). We will describe the improvement gained by using the MOTS framework on the traditional problematic of *Query By Example* (QBE) (Section 7.6.1). Then, we will show how to go beyond this conventional approach and propose two innovative audio querying methods. These new paradigms can provide easier and more intuitive control over query specifications and at the same time bypass the need for a well-formed example.

7.5.1 Content-based audio retrieval

The past decade has witnessed a growing interest in *content-based retrieval* for multimedia databases [421]. Large amount of work has been devoted to performing intuitive queries over musical *songs* databases [79], such as the *Query By Humming* (QBH) approach [433], which is now a popular content-based music retrieval method. This paradigm allows finding a song in a large collection without knowing its name or artist, just by humming its melody. Tracing back to the seminal work of Ghias et al. [143], QBH systems typically rely on *symbolic* representations of melodies, rather than generic audio databases. Sound sample databases induce a greater challenge, as they are more massive and grow faster than musical databases. Furthermore, sound samples do not benefit from the same high-level symbolic information that can be extracted from melodies. Therefore, such sets may require an overwhelming amount of time to find a specific sample. The *Query By Example* (QBE) paradigm tries to tackle this

problem by finding audio clips similar to a given sound example based on their spectral properties. The first QBE system was proposed by Wold et. al [404] where sounds were represented by a vector of spectral features, which were then compared with the Euclidean distance. This approach has subsequently been extended using larger sets of features [395] or other spectral transforms like the Discrete Cosine Transform (DCT) [365] and wavelet transform [240]. Several indexing and learning schemes have also been investigated like Nearest Feature Line (NFL) [241], Support Vector Machine (SVM) [160] or Gaussian Mixture Model (GMM) [173]. Other studies have focused on the temporal modeling of sounds, either by using templates of temporal energy [65] or Hidden Markov Model (HMM) [426] where comparison of HMM likelihoods with the query allows to obtain a ranked list of results. Finally, a different stream of generic audio querying is *Semantic Audio Retrieval* [355] which tries to discover the relationships between semantic and acoustic spaces. This provides queries on semantic concepts rather than acoustic features. This approach was implemented with a mixture of probability experts in [354] and extended with polysemy handling [69] and semantic weighting [28].

Generic audio retrieval is facing several problems that can be outlined from previous works in this field. First of all, metadata information is clearly inadequate to enable complex and intuitive interactions. It is almost impossible to maintain consistent and expressive metadata on large datasets. Semantic retrieval tries to provide a turnaround to manual annotation but still requires an extensively annotated starting set. Furthermore, it is limited to descriptive facts and sounds clearly related to a production source. However, most of the timbre ‘qualities’ cannot be captured using semantic concepts without subjective interpretation of data. Several authors pointed out the impossibility of sharing a common language for audio property description [114, 296]. This imposes severe limitations on the scope of possible queries, restricted to a predetermined set of semantic classes. Some QBE systems use clustering before retrieval, based on the idea that search time could be reduced by comparing the query only to a relevant cluster [172, 367, 427]. However, building hierarchical classes implies that the database is created according to a particular dataset. Therefore, once the database is built, it loses flexibility and users have to adapt to this original hierarchy. Finally, as we discussed earlier, authors have pointed out the multifaceted perception [114] and the unlikelihood for a single measure of perceptual similarity of audio signals [386]. Therefore, sound retrieval systems should be flexible enough so that variable influence could be put on different sound properties during perceptual similarity evaluations [264], but yet no current audio-retrieval system seems to address these limitations.

7.5.2 Going beyond traditional query paradigms

We now show how the application of the MOTS paradigm allows to handle previously listed problems. Our approach relates to [296] where sounds with or *without a known cause* are described by looking specifically at the temporal evolution of their acoustical properties. Sound clips are considered as short-duration *units of musical creativity* [78]. In order to maintain the flexibility of the database, we avoid the clustering paradigm by deliberately not interpreting data. Therefore, no assumptions are made on spectrum types and sounds can be of any nature. In order to provide more comprehensive query conditions, we do not use semantic annotation and focus on the temporal evolutions of timbre properties which provide objective comparisons.

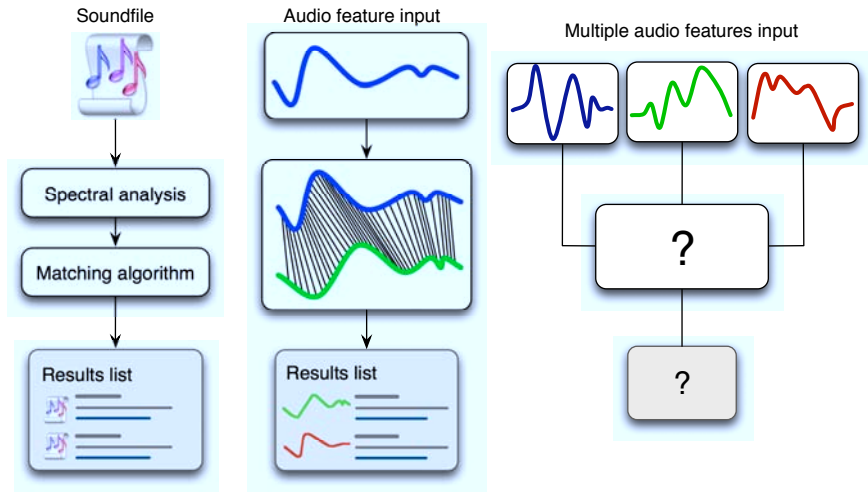


Figure 29: Shifting from the QBE paradigm (left) to the MOTS framework (right). In the QBE approach (left), a soundfile is fed to the system for which similar sounds have to be found in a database. The system answers with an ordered list of soundfile results. By using time series techniques (center), we can construct a system which match the temporal evolution of any audio feature. However, the combination of multiple audio features input (right) requires a more flexible matching process, hence exhibiting the relevance of the MOTS framework.

First, our system is based on time series data mining method (cf. Chapter 5), relying on a pre-constructed database structure (that we detail in Section 7.6). This system alone could already provide a possibility of high-level sound querying by directly matching the temporal evolution of each spectral descriptor. This would allow to find sounds by simply drawing the desired evolution of a timbre characteristic. However, as we discussed earlier, we want to confront explicitly the multidimensionality of timbre perception (cf. Section 3.3) for multiple audio features input. This paradigm shift is presented in Figure 29. In the QBE approach (left), a soundfile is fed to the system for which similar sounds have to be found in a database. The system answers with an ordered list of soundfile results. By using the time series techniques (center), we can construct a system which match the temporal evolution of any audio feature. However, the combination of several audio features input (right) requires a more flexible matching process, hence exhibiting the relevance of the MOTS framework.

Hence, we further consider a core problem of audio retrieval that lies in the query specification itself. As put forward by Donwie [114], audio queries are themselves a form of musical information, and are, therefore, complex and multifaceted. Several authors pointed out that most users of content-based retrieval systems have only a vague idea of what they seek at the onset [238, 401, 410]. Hence, they might also search for aspects of an audio query but not exactly the same content. We will show how the MOTS results handle this aspect by being presented in an informative way to users. Finally, when an example is unknown or difficult to generate, the query should help the user determine what he is seeking by being specified in a manner as close as possible to the underlying nature of audio properties [312]. To address all these shortcomings, we present two novel paradigms for audio querying based on the MOTS approach. First, the *MultiObjective Spectral Evolution Query* (MOSEQ) provides a flexible query specification by allowing users to draw directly schematic temporal

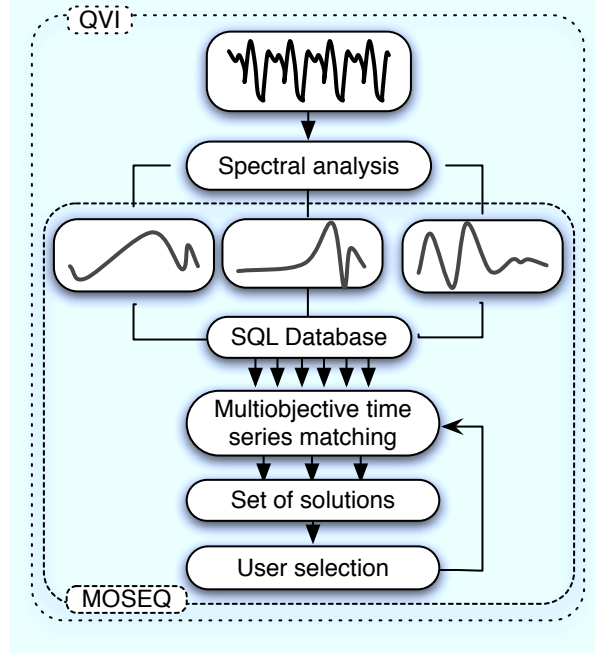


Figure 30: Algorithmic framework for two types of interaction. In *MultiObjective Spectral Evolution Query* (MOSEQ), a set of time-evolving properties is drawn. The MOTS algorithm allows to find the set of efficient solutions. In *Query by Vocal Imitation* (QVI), the user can directly use his voice to perform an imitation of the desired properties. A spectral analysis leads to the set of properties.

shapes required for spectral features. Therefore, it bypass the need for a well-formed example. Based on this paradigm, we introduce the *Query by Vocal Imitation* (QVI), which allows users to perform vocal imitations of desired properties. In both cases, the system returns the samples shown on a multidimensional front depending on how well they match the different time-evolving timbre dimensions, thus providing flexibility in results representation. Furthermore, equipped with the MOTS framework and adequate similarity measures along each perceptual dimension, we are able to predict various degrees of similarity between elements. Figure 30 summarizes the algorithmic framework for both applications.

7.6 DATABASE STRUCTURE

As we perform queries over large collections of sound samples, we have to maintain a structured database. Figure 31 depicts how sounds are analyzed and managed. We process sound samples with IRCAMDescriptor [295] in order to extract all perceptually relevant information from low-level signal data. The list of descriptors is provided in Table 2.

The mean and standard deviation of each descriptor are extracted and stored in the database. We then normalize the temporal shapes in order to obtain *zero-mean* and *unit-variance* time series. We then store the entire time series, using the SAX representation [249]. Therefore, each element in the database contains several time series which represent different characteristics of a sound. The temporal shapes are resampled to a uniform length. This could be considered as a concern for audio

Category	Features
Energy	<i>EnergyEnvelope, HarmonicEnergy, Loudness, NoiseEnergy, TotalEnergy</i>
Spectral	<i>FundamentalFrequency, Inharmonicity, Noisiness, Sharpness, Spread, Flatness, Crest, Centroid, Skewness, Kurtosis, Slope, Decrease, RollOff, Variation</i>
Harmonic	<i>Deviation, OddToEvenRatio, Tristimulus, Centroid, Spread, Skewness, Kurtosis, Slope, Decrease, RollOff, Variation</i>
Perceptual	<i>Roughness, Deviation, OddToEvenRatio, Tristimulus, Centroid, Spread, Skewness, Kurtosis, Slope, Decrease, RollOff, Variation</i>
Sub-bands	<i>MFCC, RelativeSpecificLoudness, AutoCorrelation, Chroma, ZeroCrossingRate</i>

Table 2: List of available descriptors whose mean, deviation, temporal shape and first and second derivatives are stored separately. More detailed information can be found in [295]

querying as long sounds can be compared to extremely short sounds. However, this approach shows the benefit to focus solely on the temporal shape. Furthermore, the system allows using duration in conjunction with other objectives to be optimized. The length can alternately be defined as a filtering constraint which will reduce the search space to sounds of matching length. Other symbolic information can be added to the database, either by automatic extraction from filenames or direct user input. However, we consider in the final search problem that no metadata is available whatsoever.

7.6.1 QBE results and representation

As our approach is multiobjective, query results are presented as a Pareto front in a multidimensional space. Figure 32 present the results of two queries with the MuscleFish dataset (we detail this dataset in Section 13.1.2, as it will be used to validate the MOTS framework on classification tasks). The first query (left) is performed using a restaurant scene belonging to the *crowds* class. The second query (right) is performed using a sample of *female speech*. For each query, we compare the results of mono-objective and multiobjective methods given the same set of features. Mono-objective selection provides an ordered list of results. However, there is no informed knowledge about how these choices were made whatsoever. Even with multiple dimensions involved, the results only offer an “*optimization line*” of fitness. Oppositely, the multiobjective framework allows to obtain the complete optimization space. This representation informs the user on how solutions optimize various objectives. It also allows users to explore this space by focusing more on one objective than the other. An obvious limitation of this system is on the number of dimensions that can be used for representation. However, we can easily represent three-dimensional cuts of any multidimensional space. If we look more closely at the results of these queries, we can see that the sets provided by the MOTS approach are more similar to the initial example query. In the first case, it appears that relevant results are spread over the criteria space. This distribution is revealed by the multiobjective matching. On the other hand, mono-objective selection seems to get stuck on solutions performing averagely in both objectives. Furthermore, by separating every optimization dimension, the MOTS representation already entails the case where users seeks parts of the query but not exactly the same content with

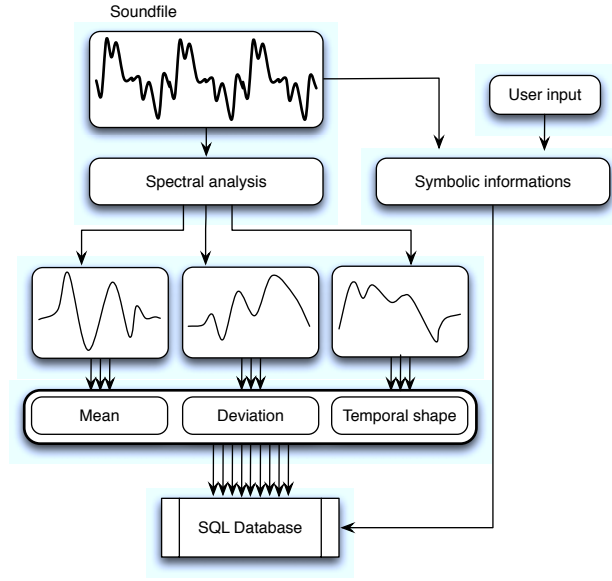


Figure 31: When a soundfile is input to the system, the analysis module computes a set of descriptors whose mean, deviation and temporal shape are stored separately inside an SQL database. Symbolic information can also be stored in the database, either by automatic extraction or direct user input.

multiple definitions of similarity. Even if this first evaluation remains an extremely narrow and empiric trial, we will validate the MOTS approach for audio querying through extensive user studies in Section 8.4.3 and complete classification tasks in Section 13.1.

7.6.2 MultiObjective Spectral Evolution Query (MOSEQ)

We present our first application of the MOTS framework that provide a novel way to find sounds in massive databases. The MOSEQ paradigm allows users to draw directly the temporal evolution of several audio features to be found in the database. Therefore, it bypass the necessity of a well-formed example and also provide a flexible query specification. We start by introducing the definitions required for this application.

Definitions

Definition 28. *Sound attributes.* An attribute \mathcal{A} is a symbolic value representing meta-informations about a sound in the database. For instance, if \mathcal{A}_{dyn} is the dynamics attribute, then $\mathcal{A}_{\text{dyn}}(\mathcal{S}) \in \{\text{pp}, \text{p}, \text{mf}, \text{f}, \text{ff}\}$. Sound attributes are obtained by manually tagging the sample database.

These may be used to define symbolic constraints on the samples and thus restrain the size of the search space. However, we do not consider symbolic attributes in the final search problem.

Definition 29. *Sound features.* A feature \mathcal{F} is a numerical (eventually multidimensional) value that describe perceptual aspects of samples. Sound features are extracted from the raw signal data (cf. Section 7.6) and can represent temporal evolutions or mean descriptors. In the latter case, they are used as attributes for reducing the search space.

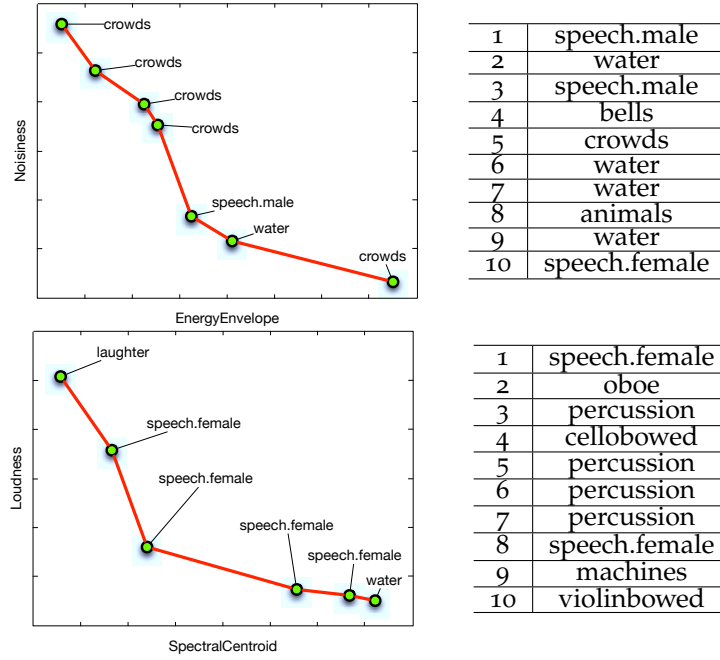


Figure 32: Comparison of different query results for multiobjective optimization and mono-objective selection in a QBE context. (Left) A sound taken from a restaurant scene and belonging to the *crowd* class. (Right) A clip taken from the *female speech* class.

We will thereafter consider the computer scientist view on timbre (cf. Section 2.3.3) and admit that any timbre may be fully characterized by a set of sound features $\{\mathcal{F}^1, \dots, \mathcal{F}^K\}$.

Definition 30. *Target.* The target query \mathcal{T} is the timbre that is input to the system and for which similar instances are to be found in the database.

For the MOSEQ system, the target is, therefore, represented by a set of time series features $\{\mathcal{F}^1(\mathcal{T}), \dots, \mathcal{F}^K(\mathcal{T})\}$ which can be of variable cardinality $K \in \{1, \dots, \mathcal{N}_{\text{features}}\}$

Definition 31. *Features similarity functions.* Given a sound target \mathcal{T} represented by its feature set $\{\mathcal{F}^1(\mathcal{T}), \dots, \mathcal{F}^K(\mathcal{T})\}$ and a sound sample \mathcal{S} , the k^{th} similarity function is the real-valued function $\mathcal{D}_{\mathcal{T}}^k(\mathcal{S})$ that returns the distance between a sound \mathcal{S} and the target \mathcal{T} along the k^{th} feature. In other words, $\mathcal{D}_{\mathcal{T}}^k(\mathcal{S})$ is a similarity measure between time series features $\mathcal{F}^k(\mathcal{S})$ and $\mathcal{F}^k(\mathcal{T})$.

Paradigm

The idea behind this interaction paradigm is that when a user seeks a sound sample, he will create a mental representation of the corresponding sound based on the temporal evolution of several spectral properties. Therefore, the MOSEQ system allows to draw these shapes in order to project a mental representation into an efficient query. The user selects a set of features that are relevant to his query. For each, he can draw the desired time series. This set acts as the target for the system. Therefore, it bypass the need for a well-formed example. We consider that the database follows the structure described in Section 7.6 and contains several sound features \mathcal{F}^i for every sample. In a QBE context, the target \mathcal{T} is the sound example for which similar instances have to be

found. For the MOSEQ system, the target is represented by a set of time series features $\{\mathcal{F}^1(\mathcal{T}), \dots, \mathcal{F}^K(\mathcal{T})\}$. As noted earlier, given this target and a sound sample \mathcal{S} , each of the similarity functions $\mathcal{D}_{\mathcal{T}}^k(\mathcal{S})$ can be defined using a different function for each objective. Therefore, the goal of the MOSEQ paradigm is to optimize simultaneously the entire set of time series features sought by the user. By using the MOTS approach, the system display the multidimensional space containing audio clips that jointly optimize the sound features. We validate the efficiency of this paradigm in Section 8.4.3 through extensive user studies and show that it easily allows to project complex time-evolving sound ideas.

7.6.3 Query by Vocal Imitation (QVI)

Now equipped with the MOSEQ paradigm, we can go even further in terms of easiness of interaction. In several context, the most straightforward way to communicate a musical idea is to use one's voice. Therefore, we believe that a natural way of querying sound samples would be to directly input a vocal imitation as a query. Indeed, the vocal system can produce a wide variety of sounds. Most people have in some occasions imitated everyday sounds by using their voice and presumably tried to match the temporal evolution of acoustic properties. For musicians, the use of nonsense text singing, called *syllabbling* [368], is an effective communication language for pedagogical purposes. Even with the inherent limitations of human voices, such as our frequency range (*tessitura*), we can control several vocal disorders. The *growl* effect increases expression by producing a rough sound. The *breathy* effect allows to generate noisier sounds. We can even learn to control extremely specific sound qualities like the position of the formant frequencies, the type of phonation or the singer's formant [369]. We thus benefit from the high degree of expression of the singing voice, principally described by loudness, fundamental frequency and spectral envelope, which all vary dynamically with time. However, it is obvious that the capabilities of a human voice are inherently limited in terms of spectral features. Nevertheless, sung imitations may convey valuable information as Pressing [308] indicates: "*One important resource in designing such expressivity is to use one's own voice to sing the expression in the part. Even if the sound quality is beyond the powers of your (or perhaps anyone's) voice, its time shaping may be imitable*". Therefore, despite the voice is limited in the range of timbres it can produce, much of vocal expression can be captured, not in the absolute timbre but in the relative temporal variation of timbre. This exhibit that our intuitions on temporal evolution can be appropriate in this setting, and we will objectively prove the efficiency and usability of the QVI framework through extensive user studies in Section 8.4.3.

Therefore, the QVI problem is defined as processing a vocal imitation provided by the user. This imitation is input to the same analysis module used for filling the database. Therefore, it provides the spectral shapes within which the user can select its desired criteria. Hence, the QVI problem can be reduced to a MOSEQ problem but still provides an equivalent of QBH for sound samples. The MOTS framework is then used to optimize the various features selected by the user jointly. We now show how to combine both paradigms and provide turnarounds to support the lack of vocal control capabilities over spectral properties.

Combining paradigms

Despite the unavoidable limitations of one's voice, it is still possible to go beyond these with two possible turnarounds. First, missing dimensions can be specified by manually

drawing the desired curves of evolution. The problem then turns to be a MOSEQ which is partly defined manually. Second, a standard mechanism used in the interaction field is to convert a control signal into another parameter (a procedure known as *mapping*). Vocal range can easily be mapped to any target range by transposing the incoming voice pitch. Hence, it is also possible to establish a mapping between any useful vocal descriptor and unrelated spectral dimensions.

We now try to validate the MOTS framework through perceptual studies and the MOSEQ and QVI paradigms through user studies. The previously presented works strongly rely on the premise that we are able to perceive and create mental representations of high-level audio features. However, even if the last decade of research have established a wealth of knowledge on the main properties of audio signals (namely *loudness* and *fundamental frequency*) the perception of more complex audio features and their temporal evolution has yet to be investigated. Furthermore, no study has tried to gain an insight on the individuality of structural organization of the multidimensional auditory perception space. We will therefore focus on studying the multidimensional perception of the temporal evolution of higher-level sound features. We start by trying to see how well temporal variations of these features can be perceived. We then study the relative perceptual importance of the features when compared to each other in a setting where they are maximally decorrelated. By relying on the concept of *directions of listening*, we try to access to the structure of our multi-dimensional space of perception. We study the consistency of these directions of listening through several tasks of direct similarity, generic similarity and shape drawing. We will show that this multidimensionality is unique to each person, even if correlations can be found in groups of subjects. We show that this multidimensional perception highly depends on the combination of features at hand. Finally, we take advantage of this multidimensional framework to perform an extensive usability evaluation of the two innovative audio querying paradigms that we introduced in the previous section. In this context, we must face experimental constraints coming from both fields of user interface (UI) and audio-related evaluations which induce large levels of variability. Therefore, we decided to conduct this study by using the *Usability Evaluation* (UE) Ellis and Dix [122] framework, which allow us to draw knowledge from the researches conducted in UE over the past years. *Usability* is the extent to which a system enables users, to achieve specified goals *effectively* and *efficiently* while promoting feelings of *satisfaction* (ISO 9421-11 [1]). We can see that this definition applies especially well in the context of querying systems. The notion of usability can be divided into three main aspects that should be investigated

- *Effectiveness* is the accuracy and completeness with which users achieve certain goals.
- *Efficiency* is the relation between effectiveness and the resources expended in achieving it.
- *Satisfaction* is the comfort and positive attitude of subjects towards the use of the system.

Several authors pointed out that usability should be evaluated on these three aspects independently Frøkjær et al. [136], an advice that we will therefore apply throughout the study.

8.1 AUDIO FEATURES

As we discussed in Section 2.3.3, a tremendous amount of work has been devoted to the development of high-level audio features. However, many of these features have been used for extremely specific task and most of these are far beyond our perceptual capabilities. For instance, the *Odd-to-Even Ratio* computes the ratio between energies of odd and even-numbered harmonic peaks. This measure, even if extremely useful for instrumental samples classification, seems to be a hard concept to grasp perceptually. Furthermore, there is a great extent of correlation between features that have been developed over the years. For instance, the *Energy Envelope*, *Loudness* and *RMS Energy* are energetic features which offer slightly different representations of the same information. Each of these features have been designed with specific goals in mind and put emphasis on certain aspects of the information. Even if they exhibit differences in computation, it is obvious that they should be strongly correlated, mostly if we focus on their temporal shapes.

As our study is centered on the perception of the temporal shapes of features, these questions are of prime importance. Furthermore, as the number of available features is way beyond the capacity of a single study, we start by selecting a subset of the most representative features that could be relevant for our purpose. We drive from these requirements, a preliminary study of correlations between features. We construct our study on the features computed with the IRCAMDescriptor Peeters [295] system. This module allows to obtain a basis of 45 temporal features (we intentionally omit multi-bands features and the various temporal modelisations). This set (listed in Table 3) can be coarsely divided into five main categories, ie. *energy*, *frequency*, *harmonic* (computed only on harmonic peaks), *spectral* (computed on the whole distribution), *noise* (computed after removing the harmonic peaks) and *perceptual* (computed after filtering the signal with a model of the human ear) features. However, these categories only provide a classification based on the computation workflow, which does not reflect the potential correlations between features which are the main interest of this study.

We start by computing all temporal features of the dataset that will be used in the perceptual study (Section 8.3.3). This dataset is composed of 3.214 environmental and synthesis sounds. Sound files are single and double channels, WAVE and AIFF format, quantized to a minimum resolution of 16-bit with a minimum sampling rate of 44.1 kHz. Loudness levels and file lengths vary with the average size of a file being about 586 Koctets. Then, for each feature, we concatenate the time series representing this feature for all the sounds in the dataset. Therefore, we obtain for each feature a long ordered time series which represent the entire sound set. The idea of concatenating the time series allow us to abstract from the duration of individual sounds when computing the correlation values. We then compute the cross-correlation between every available pairs of sound features. Therefore, we obtain 990 measurements of cross-correlations between features. Finally, we perform a hierarchical clustering on these correlations by computing the cityblock (L_1) distance and use the shortest distance between pairs to regroup features in clusters. The results of the correlation analysis are shown in Figure 33.

We see that the features exhibit strong correlations and can be organized in fewer information groups based on their correlation. We can also see that the previously cited categories does not reflect the cross-correlations between features. It seems that correlation groups are more related to the type of properties that is being computed rather than the distribution used for computation. First, we see that the most correlated group is formed by the *energetic* features. Even if putting the distinction between the

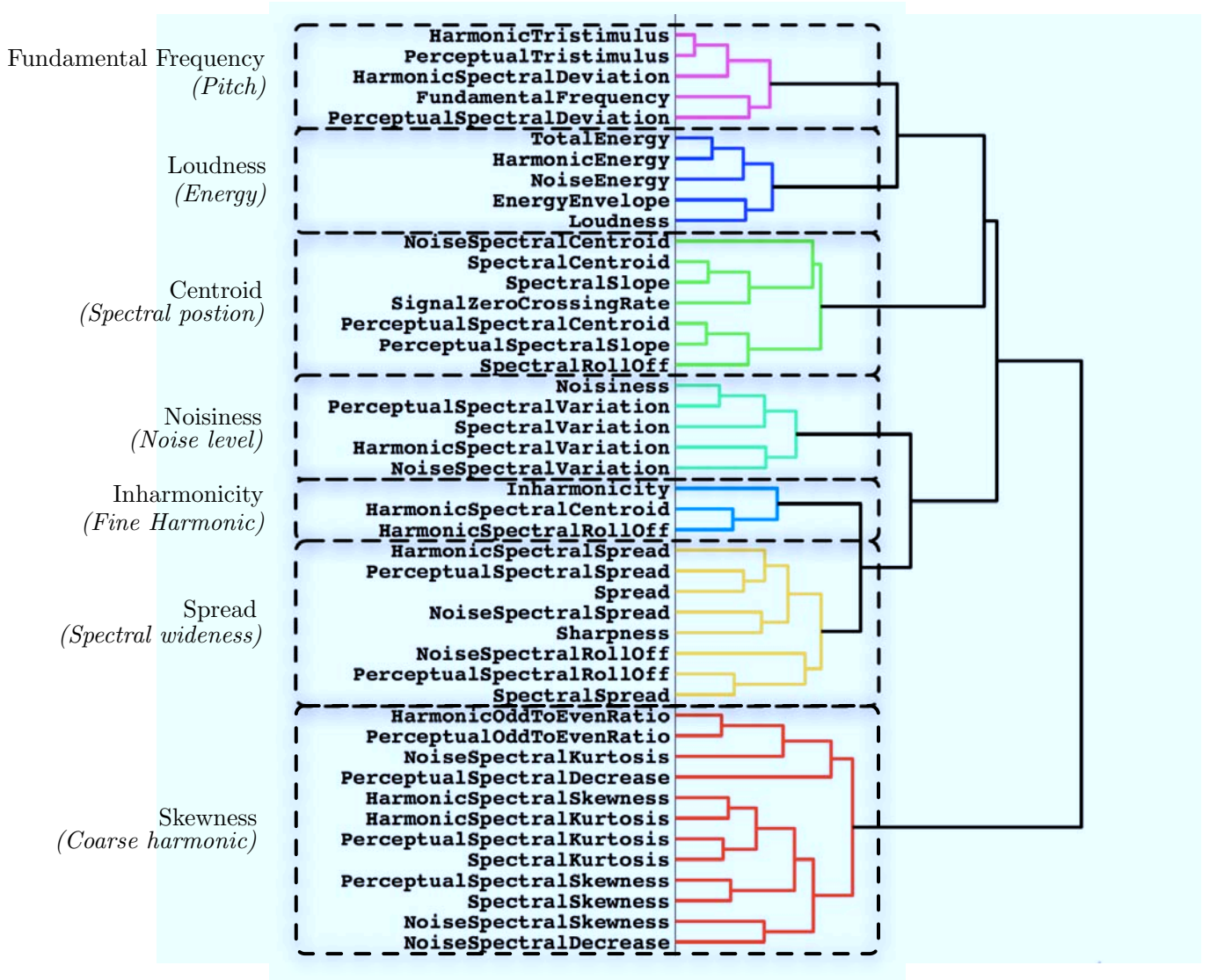


Figure 33: Cross-correlations of the temporal shapes of various audio features analyzed through a hierarchical clustering process. We represent the each cluster with its corresponding information type and representative feature selected.

energy of harmonic peaks and the energy of the noise envelope, the temporal evolution of such information appears to be strongly correlated. We therefore choose to use the *Loudness* feature to represent the *Energy* class. Another cluster that seems highly correlated relates to the information of pitch as we can see on top of the dendrogram. We choose to use the *Fundamental Frequency* feature to represent the *Pitch* class. The next group entails information on the overall position of the spectral distribution. Indeed, the *centroid*, *slope* and *roll-off* information all relate to the position of the distribution in the frequency range. We choose to use the *Spectral Centroid* feature to represent the *Position of distribution* class. Relative levels of noise in the sound are summarized by the next group which encompass the *variation* information. We choose to use the *Noisiness* feature to represent the *Levels of noise* class. The next group seems to give information on the fine harmonic structure, ie. the relative position of harmonics between each other. We choose the *Inharmonicity* feature to represent the *Fine harmonic structure* class. The two last groups seems to encompass several information regarding the overall *shape* of the spectral distribution. The first group relates to the overall width of the distribution. We choose the *Spectral Spread* feature to represent the *Width of distribution* class. Finally, the last group entails multiple information on the general structure of the spectral distribution, as summarized by the third (*skewness*) and fourth (*kurtosis*) statistical moments. We choose the *Spectral Skewness* feature to represent the *Coarse harmonic structure* class. The final selected descriptors and their corresponding correlated sets of features are provided in Table 3.

We therefore obtain a final set of seven descriptors composed of *Fundamental Frequency*, *Inharmonicity*, *Loudness*, *Noisiness*, *Spectral Centroid*, *Spectral Spread* and *Spectral Skewness* that summarize representations of different aspects of the sound information. As our study is focused on finding directions of listening *between* these different features, we will be performing each task over every possible combinations of two features, leading to a total of 21 possible pairs.

8.2 HYPOTHESES

Before starting to provide experimental protocols, it is crucial to find which questions should be raised by the evaluation framework. The first part of our study is intended at evaluating the hypotheses on which the MOTS framework was constructed. First, we want to study how different higher-level features are perceived on a temporal scale. Then, we want to assess the multidimensional nature of such temporal perception. Finally, we believe that this multidimensional organization is different for each person and the preferences over different features follow individual and consistent directions of listening.

- [Hyp.1] *Higher-level audio features can be easily perceived through their temporal evolution.*
- [Hyp.2] *The perception of the temporal evolution of different audio features follows a multidimensional organization.*
- [Hyp.3] *Preferences over spectral features are subjective, vary depending on each person and follow consistent directions of listening.*

Then, we turn our attention to the UE of our proposed audio querying paradigms. First, we divide the various interactions that constitute our approaches so that the querying system is not treated as a whole. That way, the experimental design can investigate

Information represented	Selected feature	Strongly correlated features
Pitch	<i>Fundamental Frequency</i>	Harmonic Spectral Deviation, Perceptual Spectral Deviation, Harmonic Tristimulus, Perceptual Tristimulus
Fine harmonic structure	<i>Inharmonicity</i>	Harmonic Spectral Centroid, Harmonic Spectral RollOff
Energy	<i>Loudness</i>	Energy Envelope, Harmonic Energy, Noise Energy, Total Energy
Levels of noise	<i>Noisiness</i>	Harmonic Spectral Variation, Noise Spectral Variation, Perceptual Spectral Variation, Spectral Variation
Position of distribution	<i>Spectral Centroid</i>	Noise Spectral Centroid, Perceptual Spectral Centroid, Perceptual Spectral Slope, Spectral RollOff, Spectral Slope, Signal Zero Crossing Rate
Wideness of distribution	<i>Spectral Spread</i>	Harmonic Spectral Spread, Noise Spectral Spread, Noise Spectral RollOff, Perceptual Spectral Spread, Perceptual Spectral RollOff, Sharpness, Spread
Coarse harmonic structure	<i>Spectral Skewness</i>	Harmonic Odd-To-Even Ratio, Perceptual Odd-To-Even Ratio, Harmonic Spectral Skewness, Harmonic Spectral Kurtosis, Perceptual Spectral Skewness, Perceptual Spectral Kurtosis, Noise Spectral Skewness, Noise Spectral Kurtosis, Noise Spectral Decrease, Spectral Kurtosis, Perceptual Spectral Decrease

Table 3: Selected features and information class based on the analysis of cross-correlations between sound features of the dataset used in the reminder of this study. The last column contains the corresponding set of features that are strongly correlated to the selected one.

the benefits of each of these characteristics separately. Overall, the querying systems can be divided into three main components, following our implementation workflow (Figure 30). First the nature of the *query specification*, ie. the paradigm used to input the query to the system. (in our case MOSEQ or QVI). Second, the *type of algorithm* used to find the solutions. Third, the *representation of solutions* being the full multi-objective optimization space.

- [Hyp.5] *The MOSEQ paradigm can provide an intuitive approach to project complex sound ideas into efficient queries.*
- [Hyp.6] *The QVI paradigm can provide an efficient approach to generic audio retrieval.*
- [Hyp.7] *These paradigms allow to find relevant results even without expert knowledge in signal processing.*
- [Hyp.8] *The MOTS approach allows to obtain more relevant results than linear matching.*
- [Hyp.9] *A multi-dimensional space provides a flexible representation of the proposed solutions.*
- [Hyp.10] *We are able to intuitively perform accurate vocal imitations of sound through subconscious spectral control.*

8.3 PROTOCOL

The involvement of human subjects in studies always imply several aspects of variability that can jeopardize the validity of the experiment. The study setup should thus be carefully tuned to observe the differences deriving directly from the hypotheses rather than some other uncontrolled variables. The task design also is of crucial importance to observe the system usability in response to various user behaviors. Therefore, we follow multiple guidelines provided by previous research on UE Ellis and Dix [122], Frøkjær et al. [136], Hassenzahl and Ullrich [168], Käkik and Aula [201] as well as requirements for user studies in Music Information Retrieval (MIR) Downie and Cunningham [115], Lesaffre et al. [238], Downie [114].

First, as proposed by Käkik and Aula [201], we use the *within-subjects design* where all participants perform the same set of tasks with same interfaces and same datasets, to minimize the variability induced by different users characteristics. As we are interested in studying the multidimensional perception of features, we use a *counterbalanced design* so that each features combination is used with each task equally as many times. We use *balanced feature sets*, where the number of feature sets is equal to the numbers of unique pairs that can be formed from the original set of features, leading to a total of 21 trials. To control the possible impact of the order in which these trials are presented, we counterbalance the presentation order between users using the *Latin Square* presentation.

In these tasks, we wish to obtain both *quantitative* and *qualitative* measurements, in order to take into account search performance, user satisfaction and user behavior. Therefore, we obtain *performance measurements* by processing a *log file analysis* over a *feature inspection* paradigm, where a detailed time-stamped log of all the user actions is kept and analyzed. The similarity ratings, audio played, queries performed by the users (either *vocal imitations* or *sets of drawn shapes*) as well as any mouse input are kept for further analysis. Finally, for assessing *user satisfaction*, we perform *contextual inquiries* at the end of each task by making subjects answer a small *questionnaire* for

each method. In these questionnaires, users are also able to give unconstrained *user feedback*. Hassenzahl and Ullrich [168] pointed out that retrospective evaluations *are not* consistent averages of all the experiment. Therefore, participants fill a questionnaire *after each task* to give their specific impressions on satisfaction for each method.

8.3.1 Tasks

We first try to assess the “direct” perception of higher-level audio features. Therefore, we ask users to rate the similarity between a sound heard and the temporal shapes of displayed audio features. However, we believe that this setup might introduce a *suggestive bias* by directly presenting shapes to the subjects. Therefore, we then perform a generic audio similarity task where subjects have to rate similarities between sounds without knowledge on the features used nor their temporal shapes. However, the presented sounds correspond to MOTS query results which are thus different multiobjective tradeoffs between the audio features. Finally, we want to evaluate our querying paradigms on *constrained item searching*. In these tasks, the subjects have to retrieve a specific target sound so that we can objectively measure the time required to find this sound and therefore control the variability in usability evaluation. Each task being directed as evaluating a specific hypotheses, we developed one interface specific for each task which is presented along with the task definition.

Task 1 : Shape similarity analysis

The first experiment is intended to evaluate the extent of our perception towards the temporal evolution of each feature from our selected subset but also that directions of listening may be found between features. Therefore, we want to investigate if sound features may be competing between each other when we are listening on a higher-level scale. This part of the experiment is therefore based on explicitly asking the subjects to rate the shapes of features and also to oppose two of those features at the same time, to see if one feature can take the perceptual upper hand over another.

[Task-1] *Shape similarity.* Sounds are presented with the temporal shapes of two selected audio features. Subjects are asked to rate the similarity between the sounds and the temporal shapes. Users are instructed to listen to the sounds and then rate the correlations between sounds and the temporal shapes by using maximum amplitude on a scale of [0..10] (10 indicating strong similarity).

The interface for this experiment is presented in Figure 34. In this task, for each features combination, ten sounds are presented to the user. For each of these sounds, the temporal shapes of the two features are displayed on each side of the play button. The subjects can therefore listen to each sound and rate the perceived similarity between the sound and the temporal evolution of a sound feature. Subjects are (falsely) told that the shapes may or may not correspond to the sound they heard. As the name of the features are explicitly written over each temporal evolution, the subject can focus on their shapes to rate how well he is able to hear the corresponding feature.

Task 2 : Generic similarity analysis

The second experiment has been designed to assess the perception in the selected set of sound features without any *presentation bias*. By using the multiobjective approach,

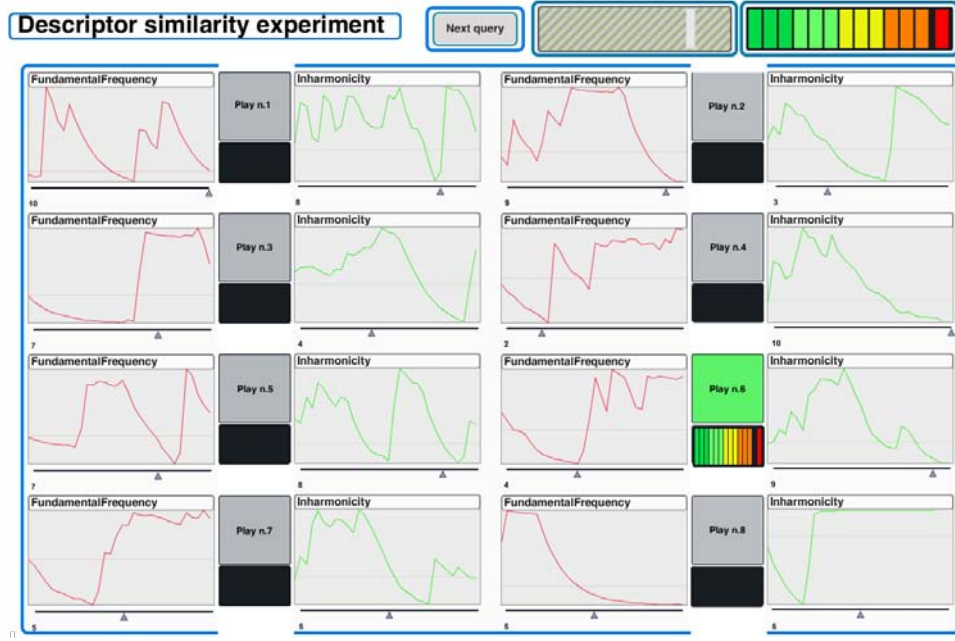


Figure 34: Experimental interface for the descriptor shape similarity task. Ten sounds are displayed for each of the 21 possible features combination. For each of these sounds, the temporal shapes of two features are displayed. Subjects are asked to rate their perceived similarity between the sound and the temporal shapes.

the MOTS algorithm is able to propose a set of sounds that represent various degrees of similarity to a query between two features simultaneously. However, for this part of the experiment, the features used and their temporal shapes are hidden to the subjects. Therefore, subjects are asked to simply rate the overall similarity between a query and a set of sounds.

[Tsk-2] *Generic similarity* : The interface presents the query and a list of sounds. Users can hear these sounds as many time as they want. Users were requested to indicate whether they agreed on the similarity with the query by rating (on a scale of [0..10] with 10 indicating that they found the sound strongly similar) the relevance of each sound proposed.

For this task, we follow the recommendations of Käki and Aula [201] and use *pre-constructed queries* and *cached result pages*. This allows to eliminate the variability induced by different setup response times and querying behavior. Therefore, we ensure that all subject can evaluate the exact same results set in a given time-frame. These problems were designed in order not to retrieve too many documents, therefore 10 potential solutions were selected for each query (by sampling the Pareto front if necessary). Therefore, for each combination of features, a query input and a list of ten sounds are presented to the subjects that they can listen as many times as they want. Subjects are asked to rate the perceived similarity between the query and the sounds by using the maximum amplitude of sliders on a scale between 0 and 10. Regarding our hypotheses, this experiment allows to assess the multidimensional hypothesis (*Hy2*) and also to study the consistency in the directions of listening (*Hy3*).

Task 3 : Constrained retrieval analysis

The last task is based on finding target sounds, known to exist inside the dataset using either the MOSEQ or QVI. In this task, the user had a sound target to attain, which he had to find back. This allows to focus on the query specification aspect, to see if a specific sound in mind is easy to find with our proposed query paradigms even for higher-level features. This task also allows to focus on evaluating the *efficiency* of the proposed query paradigms and as well to study the consistency of our directions of listening. In a experimental laboratory setup, users may take an unrealistic amount of time evaluating the result list in order to carefully follow the experimental instructions. Thus, we follow the direction of Käki and Aula [201] by providing careful instructions emphasizing that the sound might not be present in the dataset and that a similar sound (in their judgement) that can be found in the given time is acceptable. Regarding our hypotheses, these tasks allow to assess our hypotheses on the MOSEQ paradigm (*[Hyp.4]*) and the QVI paradigm (*[Hyp.5]*).

MOSEQ RETRIEVAL This part of the experiment has been designed to evaluate the usability of the MOSEQ paradigm. It has thus been constructed around a *constrained retrieval* task. In this task, subjects are presented a target sound which they have to try to find back using the MOSEQ system. Therefore, they are allowed to draw the temporal evolution of two sound features and then query the database using the MOTS algorithm. Following our previous experimental construction, subjects have to perform a query for each pair of features from our selected set. The name of the features are shown to the subjects as they are asked to draw their corresponding temporal evolution. Furthermore, compared to the first experiments, this task also offers the opportunity to see how subjects can perceive the temporal evolution of audio features *ex nihilo*.

[Tsk-2.1] *Constrained MOSEQ retrieval.* Users are provided a target sound that they try to find back by using the MOSEQ paradigm. Therefore, users are instructed to draw the temporal shapes that they think match the target on a pair of features. Each trial ends as soon as users find the target in the results list.

The interface for this experiment is presented in Figure 35. A sound target is presented to the subjects and can listened repeatedly. Subjects are asked to draw the temporal evolution of the corresponding features so that it closely match the target. The two current features are displayed under drawing boxes, which provide a breakpoint function approximation of the drawings. When subjects are satisfied with the temporal shapes, they can perform a query by using the corresponding button. The MOTS algorithm (running through OSC on a MATLAB server) will answer with a set of solutions spread over the optimization space, which provide various optimization tradeoffs to the query. Subjects can listen to the results and try to find the corresponding target. When the subjects pass over a sound result, its corresponding features are displayed. If the subjects did not find the target sound in the results, they can modify their temporal shapes and perform a new query. If the Pareto front did contain the target sound, the system displays a green success box to notify the subjects and automatically switch to the next features combination.

QVI RETRIEVAL This last task follow the same *constrained retrieval* framework to evaluate the usability of the QVI paradigm. Subjects are presented a target sound, but this time they are asked to use the QVI system to find it back. Therefore, they have to perform a *vocal imitation* of the target to query the database Subjects still have

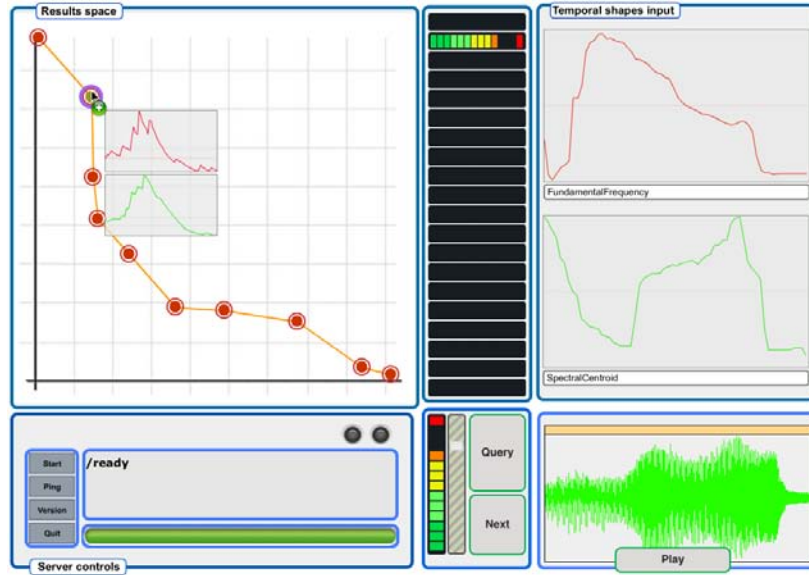


Figure 35: Interface for the constrained MOSEQ retrieval task. A sound target is presented to the subjects and can be listened repeatedly. The two current features are displayed under drawing boxes. Subjects are asked to draw the temporal evolution of the corresponding features so that it closely match the target. When subjects are satisfied with their drawings, they can perform a query. The MOTS algorithm will provide a set of solutions spread over the optimization space. Subjects can listen to the results, see the temporal evolution of their features and try to find the corresponding target.

to perform a query for each pair of features. When users perform a vocal imitation, the complete set of features analyzed from the voice recording is displayed in real-time. However, the pair of features relevant to the query are explicitly marked with a red cross. After recording their imitation, subjects are also allowed to modify the temporal shapes by drawing in the corresponding boxes. This experiment offers a unique opportunity to assess the extent of vocal control that we might possess over our selected set of audio features, based on objective comparison with the input target.

[Tsk-2.2] *Constrained QVI retrieval* : Users are now asked to find back a target sound using the QVI paradigm, i.e. they have to perform a vocal imitation. In the same fashion as before, a specific pair of features is used in each trial.

The interface for the QVI constrained retrieval task is presented in Figure 36. Overall, the behavior of this interface is the same as the previous one except for a sound analysis module that allows to display temporal shapes of sound features in real-time. Then, users can modify the shapes and perform queries in the same way as for the MOSEQ paradigm.

Experimental workflow

The experiment is therefore divided into three phases. First, each subject is introduced to the system via an informal presentation, and then asked to fill in the *introductory questionnaire*. Then the *shape similarity* interface is presented to the subjects that have to perform the similarity ratings, followed by a short break. Then the users had half an hour to perform the assessment of *generic similarity* at the end of which the *pertinency*

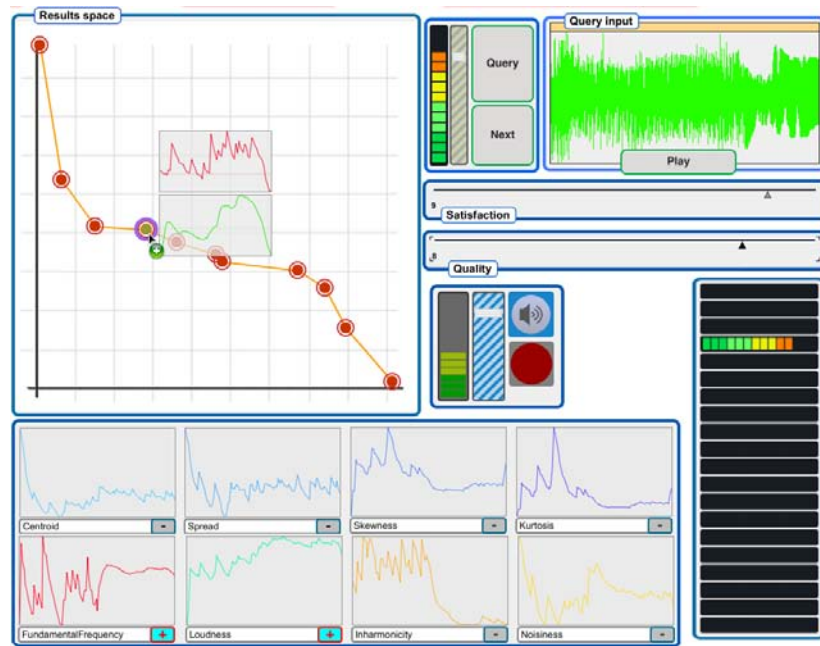


Figure 36: Interface for the constrained QVI retrieval task. A target sound is presented for each pair of features. The subjects can then perform a vocal imitation of the target. The complete set of sound features are displayed in real-time while the subject is recording. Only the pair of features relevant with the query are marked with a red cross. After recording their imitations, subjects can modify the temporal shapes by direct input. When subjects are satisfied with their input, they can perform a query. The MOTS algorithm will display the corresponding set of tradeoffs solutions.

questionnaire had to be filled in. Subsequently, the subject was given the information on the query system for the second set of tasks and could see the experimenter perform a training query which he had to reproduce. Each user had another half hour to perform these tasks and then filled in the *constrained retrieval* questionnaire after using each paradigm. This overall ordering of the experiment was controlled by an automatic activation system in sequential order through a main control panel.

8.3.2 Participants

As put forward by Ellis and Dix [122], it is important to consider the effect of the *type* of participant on the interpretation of results. Recognizing individuality is very important when trying to draw conclusions, while overall averages may hide these precious informations. Finding *to whom* a technique is really useful may be more important than making it work well for everyone and better information may be obtained with targeted groups of subjects. Furthermore, as put forward by Downie and Cunningham [115], the current MIR systems have been created with a variety of different potential users in mind but there has been no attempt yet to understand and evaluate the way that these tools get used by nonresearch communities. They stated that new research is required to better understand the specific requirements of each type of users with various backgrounds in music production. Different user groups, likely having vastly different needs, we tried to identify three types of potential users for our paradigms and formulated some hypotheses on their respective needs towards a generic audio querying system.

- *Composers* : For this group of users, the use of a generic audio querying systems may stem from its potential in stimulating creativity. Therefore the system could be used in a more flexible and exploratory manner.
- *Music producers* : For these users, audio samples are a crucial part of their workflow. Queries would therefore be aimed at finding sounds with precise characteristics. They are therefore likely to focus on the quality of the algorithm results.
- *Non-experts* : We regroup in this class of users all subjects that could find an interest in audio querying but without any formal background in musicology or signal processing. They may use the system for finding phone rings or sounds for computer actions. They are likely to focus on the simplicity of a system.

Experiments were conducted with [??] subjects ([??] female, [??] male; age [??] to [??], mean [??]), all of whom are computer-literate. Each user had to perform all the previously presented tasks. [??] subjects ([??] female, [??] male; age [??] to [??], mean [??]) without professional experience participated as "*non-expert*" users. [??] participants ([??] female, [??] male; age [??] to [??], mean [??]) had professional experience with music producing softwares and retrieval systems and participated as "*music producer*" users. Finally, [??] subjects ([??] female, [??] male; age [??] to [??], mean [??]) had professional background in musicology or composition and participated as "*composer*" users. We first investigate the homogeneity of these user groups based on their skills, collected from the introduction survey. The subjects were asked to rate their own skills in *computer use*, *music production*, *musicology*, *computer music interfaces*, *spectral features* and *singing* on a five-point Likert scale. The distribution of these skills is presented in Figure 37.

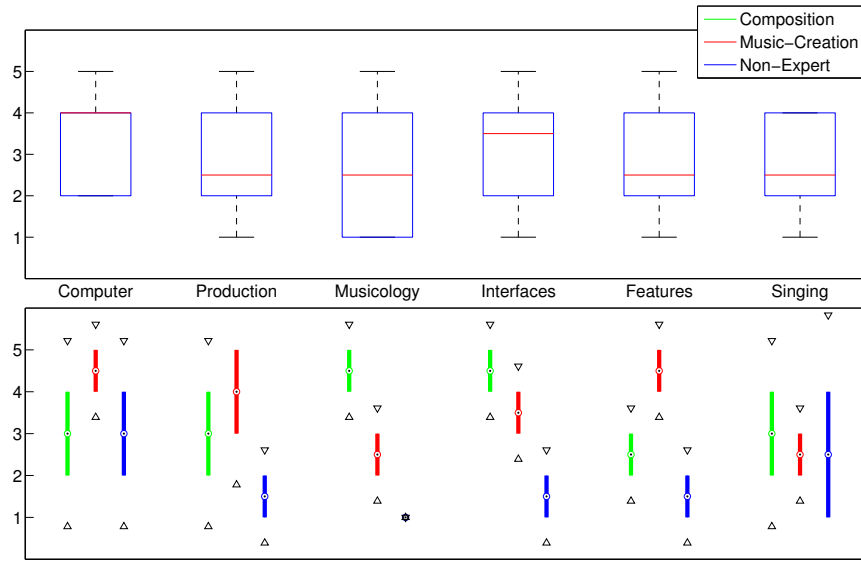


Figure 37: Distribution of skills amongst subjects of the experiment separated between mean distribution of skills (up) and distribution in each of the subjects groups (down)

As we can see, the distribution of skills is [[ANALYZE AFTER DATA COLLECTION]]

Regarding the group-wise distribution of skills, [[ANALYZE AFTER DATA COLLECTION]]

The *composition* group, [[ANALYZE AFTER DATA COLLECTION]]

The *production* group, [[ANALYZE AFTER DATA COLLECTION]]

Finally, the *non-expert* group, [[ANALYZE AFTER DATA COLLECTION]]

8.3.3 Datasets

In every audio-related experiment, the choice of the sound set is of crucial importance to truly assess the formulated hypotheses. We decided to collect a large dataset to cope with the characteristics and objectives of every tasks. The sounds were retrieved from the SoundIdea, HollywoodFX and Freesound¹ datasets. After collecting the sounds, a first filtering step was applied to ensure that the sounds were non-human non-living and could not be clearly related to a source. We want to abstract from the effects of semantic interpretation and focus on the objective temporal evolution of spectral features. After this filtering step, an automatic selection procedure was applied to further filter the dataset. First a *quality criterion* selects only high-quality sounds (quantized to a minimum resolution of 16-bit with a minimum sampling rate of 44.1 kHz). Second, a *silence detection* and removal algorithm is applied to obtain the true duration of the sound samples. As we focus on the temporal evolution, we want to abstract from widely divergent duration effects. Therefore, we keep only the sounds with durations in the range between 2 and 4 seconds. These procedures lead to a final dataset of 3,214 sounds. Sound files are single and double channels, WAVE and AIFF format, quantized to a minimum resolution of 16-bit with a minimum sampling

¹ <http://www.freesound.org>

rate of 44.1 kHz. Loudness levels and file lengths vary with the average size of a file being about 586 Koctets. As each task is intended to evaluate different hypotheses, the selection of sound queries is divided between each task.

- The first task is intended to evaluate both the perception of sound features and the multidimensional aspect of this perception when pairs of features are opposed. Therefore, for each pair of features, the 10 best sounds are automatically selected with a *maximal variance criteria* in the selected features and at the same time the *maximal decorrelation* between these features.
- For the second task, the important aspect is to focus on the multiobjective assessment of solutions. Therefore, the sound results should exhibit a wide variability. The queries are therefore selected based on a *maximal spread* of result sets over the optimization space.
- The last task focus on having the users imitate and find back a sound target. Therefore, a crucial aspect of sounds is their *reproducibility* and *distinctiveness*. We therefore apply the same automatic selection criterion as in the first task by selecting the sound with *maximal variance* in the selected features. We change the correlation criterion to select the sounds with maximal decorrelation with the rest of the set on both selected features. Furthermore, we chose to limite the cardinality of the dataset to 1.607 sounds (half of the original set) based on the maximal variance in the feature set. This choice is a tradeoff between statistical completeness and making the experiment not frustrating to the users.

8.3.4 Equipment

Each subject spent about [??] minutes in the study and received \$[??] for participation. Participants worked individually with a [??] machine with [??] DDR RAM running on [??] operating system with a high resolution [??] pixel [??]-inch color monitor, and [??] headphones in a closed isolated room in the CIRMMT laboratory at McGill university. For the QVI experiments, participants used a [??] microphone that was wired to a [??] audio interface. Participants ratings, queries, keystrokes and mouse actions were logged within the Max/Msp retrieval system during the entire session.

8.4 RESULTS

The design of each task was meant to separate our investigations over the set of hypotheses and also separate each specific aspects of usability. The analysis of results should therefore provide the same flexibility. We start by focusing the analysis of results on the perception of higher-level audio features throughout the different tasks (Section 8.4.1). After studying the overall perception of features, we further analyze the multidimensional nature of this perception. Our main assumption was that the structure of this perception is ruled by multidimensional *directions of listening* in which each subject process the same sound differently by putting variable emphasis on different features. Along this line of thought, we then further study what we call the *abstract multiobjective component* in similarity ratings (Section 8.4.2), ie. if we can find an overall correlation between the similarity ratings and the position in an optimization space. We then evaluate the usability of our new audio querying paradigms, namely MOSEQ and QVI (Section 8.4.3) by looking at their *efficiency*, *effectiveness* and *satisfaction*. Through the data collected in the usability evaluation, we provide a deeper analysis

of the querying behavior of different participants (Section ??) and their capabilities to draw the temporal evolution of audio features. The vocal imitations collected through the QVI retrieval task provide an excellent dataset to analyze the extent of spectral control in the human voice (Section 8.4.4), as exhibited by various subjects. Finally, we try to analyze the correlations and impact of the skills on these results (Section 8.4.5).

8.4.1 Multidimensional directions of listening

Our goal in the first part of the experiment was to find the structural organization of the multidimensional perceptual space for each individual. We consider that sound features might be opposed with variable influence on each other, ie. some sound features can shadow the perception of others. Therefore, we try to compute the relative *perceptual strength* between different sound features. We hypothesized that users will have different preferences between features (and therefore different structural organizations for their multi-dimensional space of perception).

The [Tsk-1] task explicitly ask the participants to rate the correlations so we can straightforwardly analyze to what extent the participants can perceive the temporal shape of a particular sound feature by using the collected similarity ratings. However, this task might exhibit an unwanted bias that arise from the mental suggestion imposed on the subjects. Indeed, the participants know that a presented set of temporal shapes could be related to the sound they hear, so they might try to force themselves to consider this correlation. The remaining task has however implicit ways of providing turnarounds to this problem and can further confirm the validity of our hypotheses. Indeed, when users are asked to rate generic sound similarity without knowledge on the features nor their temporal shapes in the [Tsk-2] task, these perceptual judgements are exempt of suggestive bias. Afterwards, when subjects have to draw temporal shapes or perform vocal imitations in [Tsk-3], they will simply draw the mental image of potential features of the sound that they just heard. Therefore we can study how users perceive the temporal evolution of spectral features *ex nihilo* (and without any suggestive bias). However, these last tasks do not provide a direct access to a measure of *perceptual strength* over sound features. Therefore, we detail the pre-processing applied to the data in order to obtain such information for each task in its corresponding part.

Shape similarity

In this task, the displayed temporal shapes are the true output of the audio features analysis module. Therefore, we can use the user ratings straightforwardly as a measure of perceptual strength.

GENERIC FEATURE PERCEPTION We start by analyzing the overall distribution of similarity ratings. In this preliminary analysis, we do not take into account *which combinations* of features are studied. Therefore, we abstract from the multi-dimensional aspect of this study. We simply process all ratings from every subjects and for each feature to get insights on the overall perception of higher-level sound features. As features were involved in 6 combinations each, with 10 sounds presented each time, we obtain a total of 60 scores per subject for each feature. We compute the mean similarity score per subject for each feature. The distribution of subjects ratings and group-wise ratings are presented in Figure 38.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

[[FEATURE-DEPENDENT ANALYSIS]]

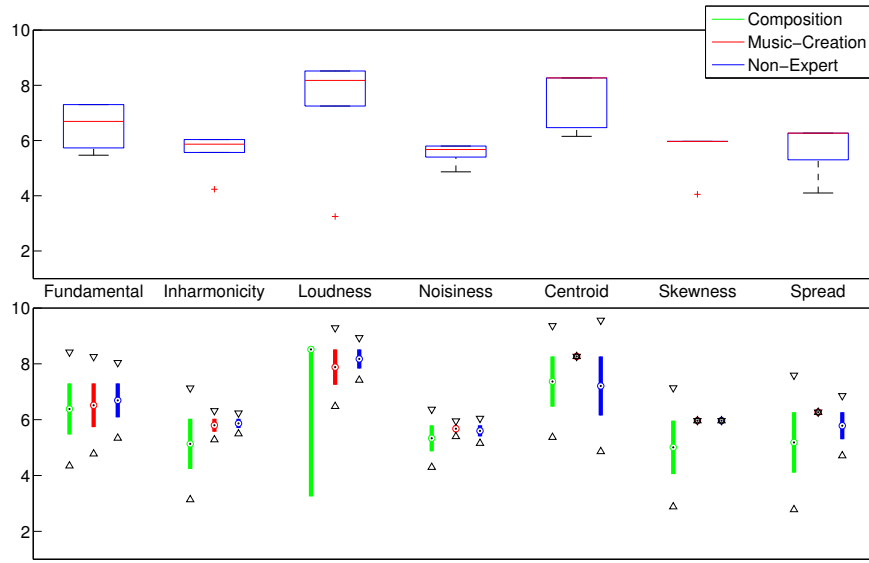


Figure 38: Distributions of mean similarity ratings for each sound feature, independently of the combinations used. Similarity scores for all subjects (up) and group-wise similarity ratings (down).

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[COMPARISON OF FEATURE SCORES]]

[[VARIANCE-BASED ANALYSIS]]

[[GROUP-WISE ANALYSIS]]

DIRECTIONS OF LISTENING We now analyze our hypothesis on the multidimensional perception of audio features through the concept of *directions of listening*. We construct this analysis around the idea that features might be competing with each other on a perceptual level, ie. perception of the temporal evolution of a feature depends on the variation of other features. For instance, the perception of *Spectral Spread* might be shadowed by variations in *Fundamental Frequency* that drifts away our attention. These directions of listening might be widely different from one person to the other. Therefore, each person might have inconscious preferences of listening when trying to process audio features. In order to exhibit this behavior, we process the similarity ratings in a tournament-based fashion, ie. we consider that directions of listening depends on the opposition between two features of a combination. We will then study the consistency and tasks effects of these directions throughout the study.

We start by trying to find the overall directions of listening, to see if a global ordering can emerge from the ratings of all subjects. To that end, we start by computing the tournament matrix for all subjects. Therefore, for each subject, we compute the difference between the ratings of both features of each combination. Then, we sum these differences for the ten sounds of each combination, which allows to obtain a triangular matrix of tournament. This matrix represents the relative weights of preference between each pair of features. We fill the other part of the matrix with the corresponding opposite values and then normalize the complete matrix. We then apply

	Fund.		Inha.		Loud.		Nois.		Centro.		Skewn.		Spread	
	#		Val	#	Val	#	Val	#	Val	#	Val	#	Val	#
Similarity	6.53	3	5.60	5	6.53	3	5.60	5	6.53	3	5.60	5	5.60	5
Composers	6.38	3	5.13	5	6.38	3	5.13	5	6.38	3	5.13	5	5.13	5
Production	6.51	3	5.80	5	6.51	3	5.80	5	6.51	3	5.80	5	5.80	5
Non-Experts	6.69	3	5.87	5	6.69	3	5.87	5	6.69	3	5.87	5	5.87	5
Score	5.45	5	-7.16	3	5.45	5	-7.16	3	5.45	5	-7.16	3	-7.16	3
Composers	7.75	5	-8.4	1	7.75	5	-8.4	1	7.75	5	-8.4	1	-8.4	1
Production	1.65	5	-8.1	3	1.65	5	-8.1	3	1.65	5	-8.1	3	-8.1	3
Non-Experts	6.95	3	-5.0	5	6.95	3	-5.0	5	6.95	3	-5.0	5	-5.0	5
Eigenvector	0.0	1	-0.0	6	0.0	1	-0.0	6	0.0	1	-0.0	6	-0.0	6
Composers	0.32	1	0.32	3	0.32	1	0.32	3	0.32	1	0.32	3	0.32	3
Production	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	3
Non-Experts	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2
Variance	0.94	1	1.56	2	0.94	1	1.56	2	0.94	1	1.56	2	1.56	2
Composers	0.94	1	1.88	2	0.94	1	1.88	2	0.94	1	1.88	2	1.88	2
Production	0.94	1	1.58	2	0.94	1	1.58	2	0.94	1	1.58	2	1.58	2
Non-Experts	0.94	1	1.52	2	0.94	1	1.52	2	0.94	1	1.52	2	1.52	2

Table 4: Tournament-based analysis of listening directions for the *shape similarity* task.

three different procedures on these tournament matrix as presented in Table 4 for the shape similarity task.

First, we sum up the columns of the matrix and take the mean vector of all subjects. This analysis is the *score-based* procedure, in which the features are ordered depending on the overall preferences exhibited through every combinations. This procedure gives us the overall weight and rank of each feature against every other, ie. which feature has been consistently rated superior or inferior to the other features. In this case, the sum of a row (score value) can be interpreted as the *expected similarity difference* for a feature against any other. Second, we compute the eigenvectors of the tournament matrix. This analysis gives the *eigenvector-based* procedure, in which features are ordered depending on their corresponding components in the eigenvector. In this case, the Perron-Frobenius theorem guarantees the existence of an unique large eigenvalue and that the eigenvector has strictly positive components. The component corresponding to a feature in the eigenvector represents a measure of its relative *strength*, as shown by Keener [206]. Finally, for each feature we take the variance of its relative scores when compared to other features. This analysis gives the *variance-based* procedure, in which features are ordered depending on the variance of their scores. This procedure is intended to exhibit which features are least affected by other features, ie. if their similarity ratings are independent of the combinations studied.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

[[SCORE-BASED ANALYSIS]]

[[EIGENVECTOR-BASED ANALYSIS]]

[[VARIANCE-BASED ANALYSIS]]

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]

In order to obtain a detailed view over the *across-subjects* distribution of directions of listening, we now consider each pair of ratings from every subjects as a single competition. We compute the scatter plots and kernel density estimates of every similarity ratings for all features combinations and represent these results in Figure 39. This figure exhibits the complete matrix of pairwise distributions of all similarity ratings. The lower triangular matrix displays the scatter plots of every ratings, depending

on the group of the subject. For each entry of this part of the matrix, the column feature is displayed as abscissa and the row feature is displayed as ordinates. For instance, the top left entry of the scatter plots displays the rating scores of *Fundamental Frequency* as abscissa against the ratings of *Loudness* as ordinates. The upper triangular matrix uses the same information with reversed coordinates but displays the kernel density estimates of the similarity ratings. These entries thus give a straightforward overview of the similarity regions where most of the ratings are concentrated, ie. on the overall directions of listening. Therefore, the lower matrix gives insights on the group-dependent distribution of ratings, whereas the upper matrix provides an overview of the concentration of ratings in every directions of listening, ie. how features are competing to each other in a pairwise fashion.

```

[[ NEED TO
[[ MODIFY
[[ THIS FIGURE
[[ : Encadrer les features ambiguës
[[ : Encadrer les pas du tout ambiguës
[[ : Encadrer les all-winners (flèche pour montrer la disparité)
[[ : Encadrer les all-losers
[[ ANALYSIS OF OVERALL DISTRIBUTIONS ]]
[[ COMPARISON OF FEATURE SCORES : Compare very narrow spots (com-
pletely oriented listening) vs. very wide-spread distributed kernels (high variance
of listening directions) ]]
[[ DOGMATIC FEATURES : Features that win all the time (Loudness ?) ]]
[[ LITIGIOUS FEATURES : Features that balance each other ]]
[[ LAME FEATURES : Features that loose ]]
[[ VARIANCE-BASED ANALYSIS ]]
[[ FEATURE-DEPENDENT ANALYSIS ]]
Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spec-
tral Skewness - Spectral Spread
[[ GROUP-WISE ANALYSIS ]]

```

The previous analysis allowed us to gain a global overview on the distributions of pairwise similarity ratings and therefore the relative weights of features between each other. However, with this representation we can not truly discriminate the individual differences that may arise in perceptual judgements. We now focus on the second part of the multidimensional hypothesis being that listening preferences should be mostly unique to each individual. We provide the scatter plots and kernel density estimates from an individual analysis in Figure 40. The display has the same organization as previously, however the input is slightly different. We use the normalized mean differences between similarity scores for each subject. That way, each point of the scatter plots represent the relative listening directions of a particular subject towards each pair of features.

```

[[ ANALYSIS OF OVERALL DISTRIBUTIONS ]]
[[ COMPARISON OF FEATURE SCORES : Compare very narrow spots (com-
pletely oriented listening) vs. very wide-spread distributed kernels (high variance
of listening directions) ]]
[[ DOGMATIC FEATURES : Features that win all the time (Loudness ?) ]]
[[ LITIGIOUS FEATURES : Features that balance each other ]]
[[ LAME FEATURES : Features that loose ]]
[[ VARIANCE-BASED ANALYSIS ]]
[[ FEATURE-DEPENDENT ANALYSIS ]]

```

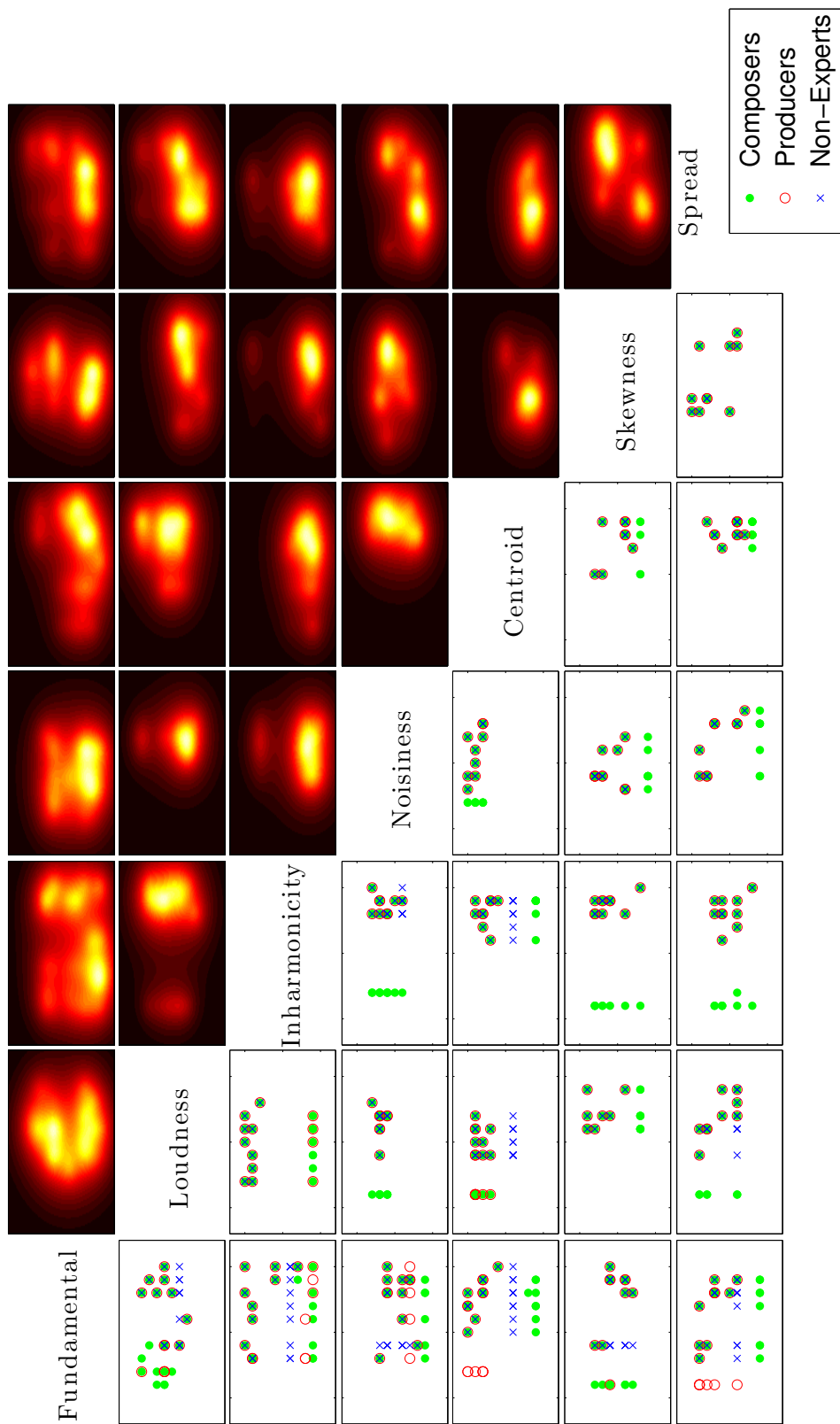


Figure 39: Scatter plots and kernel density estimates for all similarity ratings available and each feature combination.

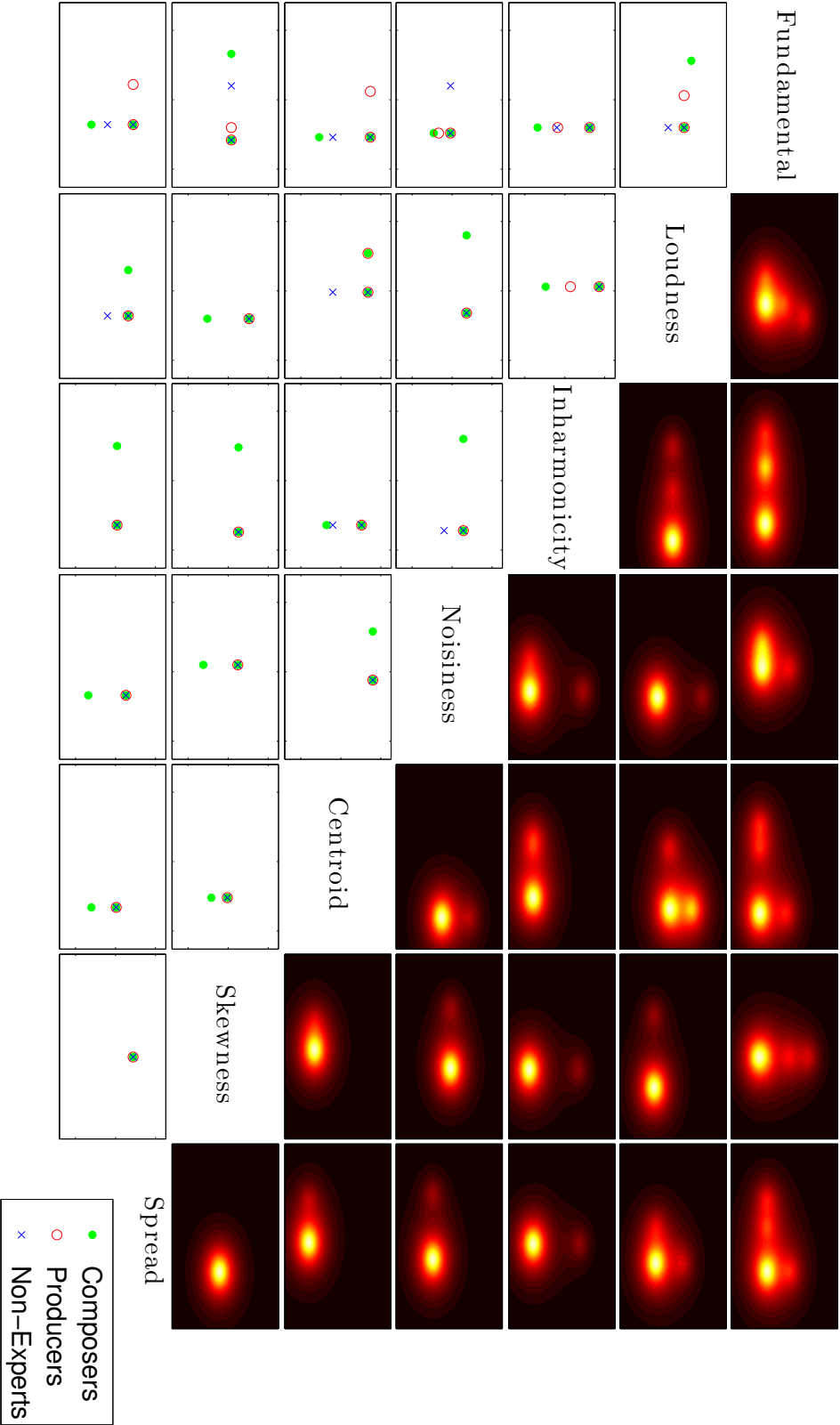


Figure 40: Scatter plots and kernel density estimates for mean similarity ratings produced by each subject.

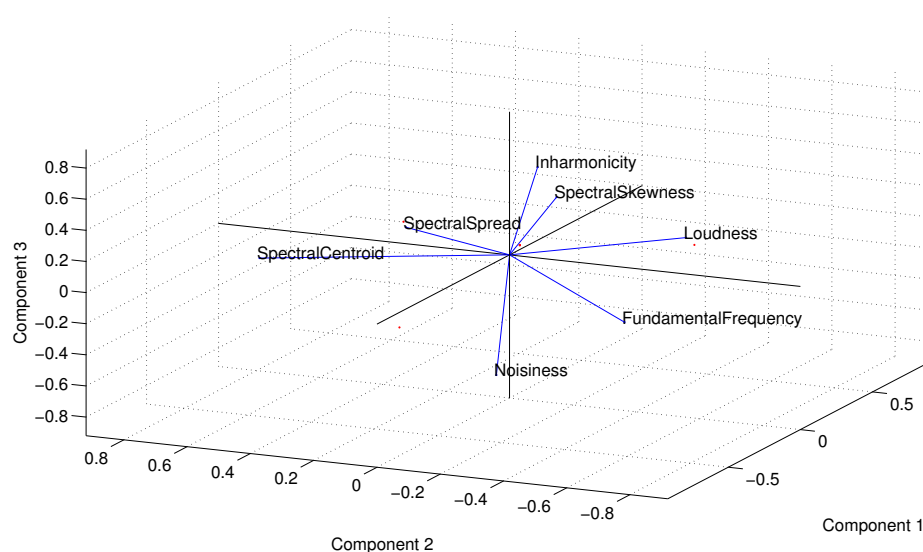


Figure 41: Principal Components Analysis (PCA) of the similarity tournament matrix.

	Var.	Fund.	Loud.	Inha.	Nois.	Cent.	Spre.	Skew.
1								
2								
3								

Table 5: Linear coefficients of the PCA for the three main components displayed in Figure 41

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]

Given the previous analyses on the individuality of directions of listening, we now try to extract the main axis that could summarize the variability found in the similarity ratings. This type of data reduction could lead to explain some prototypical *strategies of listening*. We therefore apply a Principal Component Analysis (PCA) on the user similarity tournament matrix. The results of this analysis are presented in Figure 41.

[[VARIANCE EXPLAINED - OVERALL ANALYSIS]]

[[ANALYZE FIRST COMPONENT]]

[[ANALYZE SECOND COMPONENT]]

[[ANALYZE THIRD COMPONENT]]

[[PCA LINEAR COEFFICIENTS ANALYSIS]]: We provide in Table 5 the linear coefficients for the three main components of the PCA displayed in Figure 41. This allows to see the precise contribution of each feature to the principal components that explains the prototypical directions of listening.

[[PROTOTYPICAL DIRECTIONS OF LISTENING]]

[[USER DISTRIBUTION ANALYSIS]]

[[FEATURE-WISE ANALYSIS]]

Now that we have gained a global view on the structural organization and individual differences in the directions of listening, we would like to see if we can extract trends or

similar behaviors between subjects and groups. These patterns of similarities can help us identify the perception of high-level features. However, the data we need to analyze is already multivariate for each subject and we wish to extract potential correlations between several subjects and groups for multiple features combinations. Therefore, we use a redundancy analysis (RDA) to examine such correlations. This technique performs a constrained ordination that exhibits how the variation in one set of variables explains the variation in another set of variables. It is the multivariate analog of linear regression based on a PCA. We apply the RDA on the tournament information obtained from last steps. More precisely, for each subject we compute the complete tournament matrix for all features combinations (cf. Figure 40). We store this matrix as we will use it later to check the consistency of listening directions throughout remaining tasks. We then process each individual matrix to obtain the *score-based tournament vector* (analogous to the results presented in Table 4). The final data matrix thus contains for each subject a vector of weights which represent its individual directions of listening. Therefore, each subject is treated as an *observation* and each sound feature as a *variable*. We first perform the RDA by considering the user groups categorization. The results are presented in Figure 42.

To compute the RDA, we use the Fathom toolbox Jones [198]. In this implementation, eigenvectors are scaled to unitary lengths in order to preserve the distances among subjects. The figures are to be interpreted as follows. Distances among subjects approximate their Euclidean distance. Projecting subjects onto a feature arrow approximates the mean similarity rating of that subject along the corresponding feature. Angles between groups and features reflects their correlations. We obtain a F-Statistic of [????] with p-value of [????] and coefficient of determination $R^2 = [????]$.

[[VARIANCE EXPLAINED - OVERALL ANALYSIS]]

[[ANALYZE FIRST CANONICAL AXIS]]

[[ANALYZE SECOND CANONICAL AXIS]]

[[GROUP-WISE ANALYSIS]]

[[OUTLIERS ANALYSIS]]

[[EIGENVALUES ANALYSIS]]

[[FEATURE-WISE ANALYSIS]]

Now that we have analyzed the correlations between different groups, we try to extract the same kind of information by looking specifically at redundancies between subjects. Furthermore, this analysis allow us to confirm if our main hypothesis is valid throughout different subjects. Indeed, this analysis would exhibit some typical disparities if various listening strategies are applied by different subjects. Therefore, we perform a RDA with the same pre-processing and parameter settings. The results of this analysis is presented in Figure 43.

We obtain a F-Statistic of [????] with p-value of [????] and coefficient of determination $R^2 = [????]$.

[[VARIANCE EXPLAINED - OVERALL ANALYSIS]]

[[ANALYZE FIRST CANONICAL AXIS]]

[[ANALYZE SECOND CANONICAL AXIS]]

[[GROUP-WISE ANALYSIS]]

[[OUTLIERS ANALYSIS]]

[[FEATURE-WISE ANALYSIS]]

[[CONCLUSIONS ON DIRECTIONS OF LISTENING]]

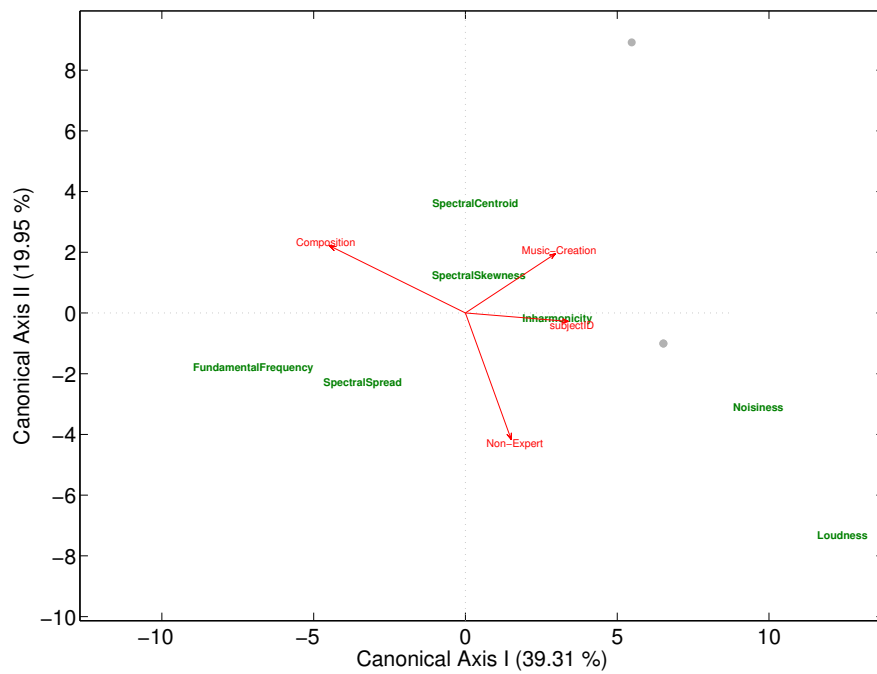
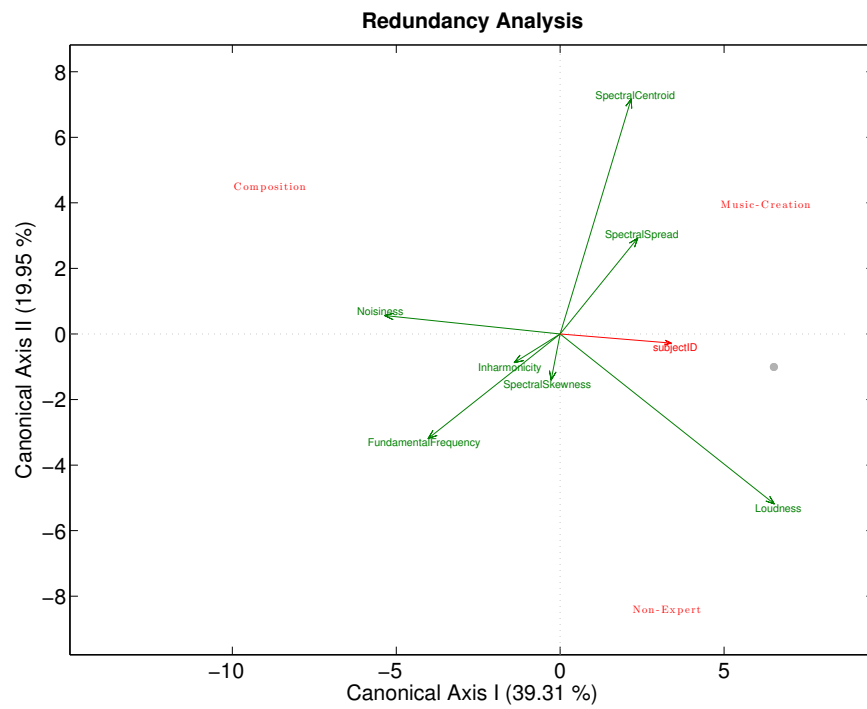


Figure 42: Redundancy Analysis (RDA) results performed over the individual score-based directions of listening depending on the groups.

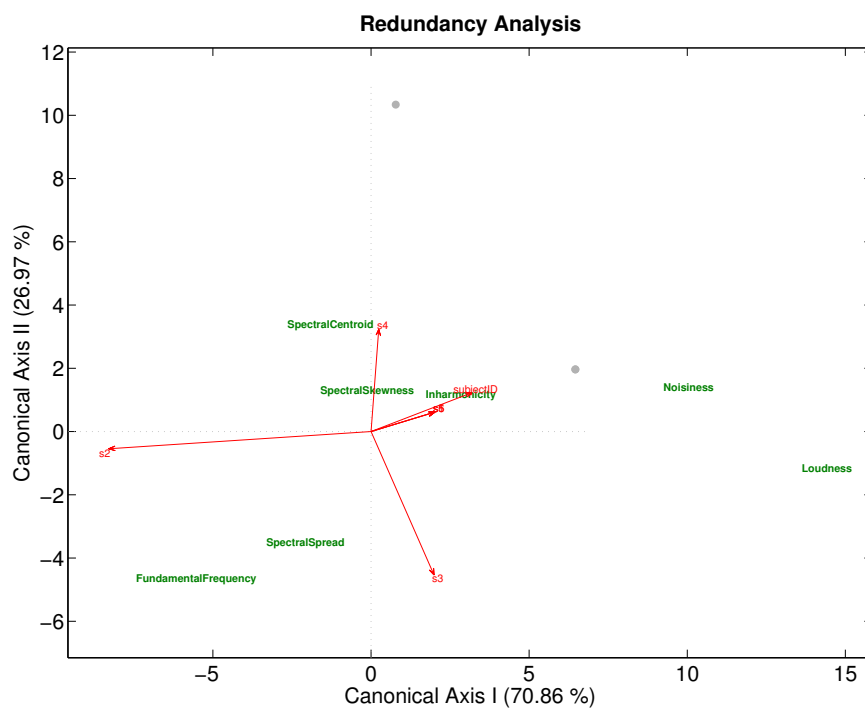
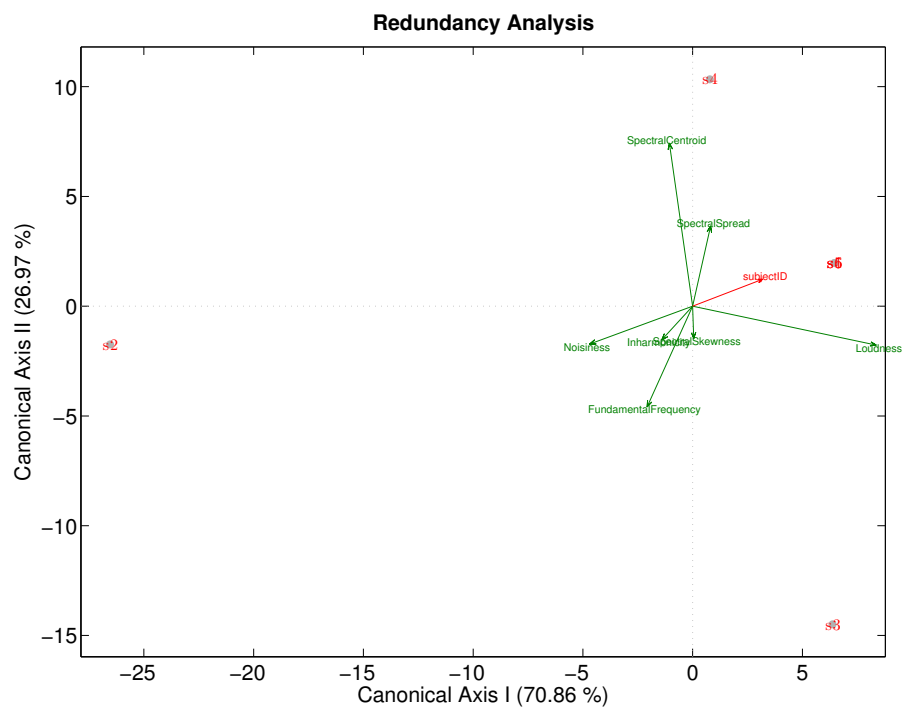


Figure 43: Redundancy Analysis (RDA) results performed over the individual score-based directions of listening depending on the subjects.

	Fund.		Inha.		Loud.		Nois.		Centro.		Skewn.		Spread	
	Val	#	Val	#	Val	#	Val	#	Val	#	Val	#	Val	#
Similarity	6.53	3	5.60	5	6.53	3	5.60	5	6.53	3	5.60	5	5.60	5
Composers	6.38	3	5.13	5	6.38	3	5.13	5	6.38	3	5.13	5	5.13	5
Music-Creation	6.51	3	5.80	5	6.51	3	5.80	5	6.51	3	5.80	5	5.80	5
Non-Experts	6.69	3	5.87	5	6.69	3	5.87	5	6.69	3	5.87	5	5.87	5
Score-based	5.45	5	-7.16	3	5.45	5	-7.16	3	5.45	5	-7.16	3	-7.16	3
Composers	7.75	5	-8.4	1	7.75	5	-8.4	1	7.75	5	-8.4	1	-8.4	1
Music-Creation	1.65	5	-8.1	3	1.65	5	-8.1	3	1.65	5	-8.1	3	-8.1	3
Non-Experts	6.95	3	-5.0	5	6.95	3	-5.0	5	6.95	3	-5.0	5	-5.0	5
Eigenvector	0.0	1	-0.0	6	0.0	1	-0.0	6	0.0	1	-0.0	6	-0.0	6
Composers	0.32	1	0.32	3	0.32	1	0.32	3	0.32	1	0.32	3	0.32	3
Music-Creation	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	3
Non-Experts	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2
Variance	0.94	1	1.56	2	0.94	1	1.56	2	0.94	1	1.56	2	1.56	2
Composers	0.94	1	1.88	2	0.94	1	1.88	2	0.94	1	1.88	2	1.88	2
Music-Creation	0.94	1	1.58	2	0.94	1	1.58	2	0.94	1	1.58	2	1.58	2
Non-Experts	0.94	1	1.52	2	0.94	1	1.52	2	0.94	1	1.52	2	1.52	2

Table 6: Tournament-based analysis of listening directions for the generic similarity task.

Consistency in generic similarity

In the second task, users were asked to rate the similarity between a query sound and a set of results without knowing which sound features were involved in the similarity computation nor their corresponding temporal shapes. We assess our main hypothesis by analyzing if the directions of listening inferred from the previous step are consistent with this task. Therefore, we will see if the directions of listening of a subject is a coherent and persistent phenomenon which still applies in generic sound similarity. We first analyze the directions of listening that can be extracted from the generic similarity ratings. We then try to correlate these directions of listening for each subject with those found in the previous task. However, the experimental setting for this task is widely different from the previous one and does not provide a direct access to the *perceptual strength* information for each feature. Therefore, we perform a pre-processing step in order to infer the directions of listening from the similarity ratings. The main idea behind this extraction method is that each sound in the Pareto front have a position representing its distance measures in each objective. We can compute the normalized direction vector between the origin of the space and a sound solution. This sound would be the best solution of a mono-objective problem where the dimensions are weighted by the components of this vector. Therefore, we consider the *position-weighted similarity score* as a measure of perceptual strength. In this method, we consider that each similarity rating applies to both features of the combination (dimensions of the corresponding space), weighed by its position vector. Therefore, for each feature, we sum the similarity ratings weighted by the normalized positions in their corresponding dimension. As before, we first try to find the overall directions of listening by computing the tournament matrix for all subjects. We then apply the same *score-based*, *eigenvector-based* and *variance-based* procedures on the tournament matrix. The results of this analysis are presented in Table 6.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

[[SCORE-BASED ANALYSIS]]

[[EIGENVECTOR-BASED ANALYSIS]]

[[VARIANCE-BASED ANALYSIS]]

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]

We now focus our attention on the individual directions of listening. We apply the same computational steps as previously in order to obtain the individual tournament matrix for each subject. Our goal is now to analyze the consistency between directions of listening when shifting from a “direct” rating of shape similarity to the generic sound similarity. Therefore, we apply the RDA method by still considering the subjects of the experiments as *observations* and the different features as *variables*. However, this time we consider the previous set of listening preferences to be the *explanatory variable*. Therefore, the algorithm will look for redundancies between the two sets of listening directions, which will allow us to see if the ratings in each feature are correlated. It will therefore exhibit if the directions of listening are coherent and consistent between the tasks. We present the results of this analysis in Figure 44.

We obtain a F-Statistic of [????] with p-value of [????] and coefficient of determination $R^2 = [????]$.

[[VARIANCE EXPLAINED - OVERALL ANALYSIS]]

[[FEATURES CORRELATION ANALYSIS]]

[[FEATURE-WISE ANALYSIS]]

[[ANALYZE FIRST CANONICAL AXIS]]

[[ANALYZE SECOND CANONICAL AXIS]]

[[GROUP-WISE ANALYSIS]]

Constrained retrieval

Even if the main goal of the constrained retrieval tasks is to evaluate the usability of the audio querying paradigms, they also provide a wealthy source of information for perceptual analysis. Indeed, these tasks offer an opportunity to analyze the perception of sound features from another angle, by inferring the *perceptual strength* from the drawn queries. As the previous task was intended to abstract from the feature information, this task now directly assess how subjects perceive the temporal evolution of sound features *ex nihilo*. Therefore, we can see how similar the drawn shapes are to the original sound features. However, like the previous task, we need to find a way to extract these measurements. This can be done by considering the time series distance as a measure of perceptual dissimilarity. We therefore study the distances between the shapes drawn by the users and the true sound features computed over the queries, by using the *Dynamic Time Warping (DTW)* distance. As two features are required for each query, we also apply these measures in a tournament fashion (as performed previously). Therefore, after computing the distance between the true sound features and the drawn shapes, we normalize these distances and then apply the same computation steps as previously. The corresponding results are presented in Table 7.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

[[SCORE-BASED ANALYSIS]]

[[EIGENVECTOR-BASED ANALYSIS]]

[[VARIANCE-BASED ANALYSIS]]

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]

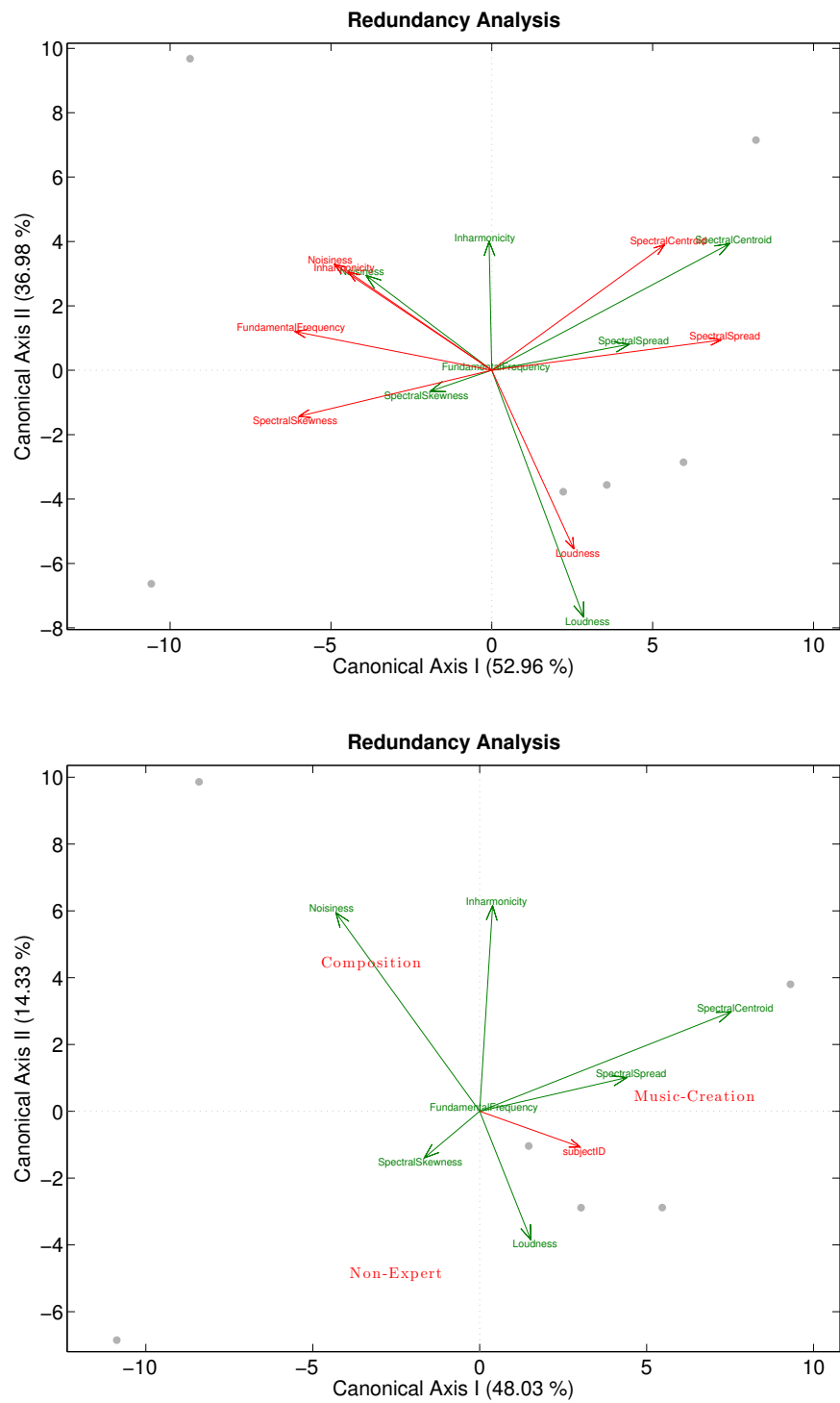


Figure 44: Analysis of the consistency in the directions of listening using a RDA method for the generic similarity task.

	Fund.		Inha.		Loud.		Nois.		Centro.		Skewn.		Spread	
	Val	#	Val	#	Val	#	Val	#	Val	#	Val	#	Val	#
Similarity	6.53	3	5.60	5	6.53	3	5.60	5	6.53	3	5.60	5	5.60	5
Composers	6.38	3	5.13	5	6.38	3	5.13	5	6.38	3	5.13	5	5.13	5
Music-Creation	6.51	3	5.80	5	6.51	3	5.80	5	6.51	3	5.80	5	5.80	5
Non-Experts	6.69	3	5.87	5	6.69	3	5.87	5	6.69	3	5.87	5	5.87	5
Score	5.45	5	-7.16	3	5.45	5	-7.16	3	5.45	5	-7.16	3	-7.16	3
Composers	7.75	5	-8.4	1	7.75	5	-8.4	1	7.75	5	-8.4	1	-8.4	1
Music-Creation	1.65	5	-8.1	3	1.65	5	-8.1	3	1.65	5	-8.1	3	-8.1	3
Non-Experts	6.95	3	-5.0	5	6.95	3	-5.0	5	6.95	3	-5.0	5	-5.0	5
Eigenvector	0.0	1	-0.0	6	0.0	1	-0.0	6	0.0	1	-0.0	6	-0.0	6
Composers	0.32	1	0.32	3	0.32	1	0.32	3	0.32	1	0.32	3	0.32	3
Music-Creation	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	3
Non-Experts	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2
Variance	0.94	1	1.56	2	0.94	1	1.56	2	0.94	1	1.56	2	1.56	2
Composers	0.94	1	1.88	2	0.94	1	1.88	2	0.94	1	1.88	2	1.88	2
Music-Creation	0.94	1	1.58	2	0.94	1	1.58	2	0.94	1	1.58	2	1.58	2
Non-Experts	0.94	1	1.52	2	0.94	1	1.52	2	0.94	1	1.52	2	1.52	2

Table 7: Tournament-based analysis of listening directions for the constrained retrieval task.

8.4.2 Abstract multi-dimensional similarity

We try to determine here if the application of multiobjective principles is relevant to sound similarity. Therefore, we want to see if the flexibility introduced by new concepts of similarity is impacted on the ratings provided by users. Indeed, the solutions presented in the generic similarity task offer different tradeoffs between dimensions. Therefore, we analyze if we can observe an overall spatial preference over the multiobjective space, independently of the features used in the combination. That is, if subjects consistently give higher similarity ratings towards equal-weighted mono-objective solutions, ie. towards the middle of the space. These solutions are the results of usual nearest-neighbor searches that performs a linear mix of the similarities in every features. Higher similarity ratings concentrated in this region of the optimization space would imply that the multiobjective framework might not be of relevance for audio similarity perception. Otherwise, we can see if users are consistently giving higher similarity to extreme solutions that optimize well only one objective or any other regions that deviate from the equally-weighted mono-objective solutions. We will see if there is an *overall correlation* between the position of the sounds in the optimization space and the similarity ratings, ie. if relevant document are generally more at the edges or in the middle of the Pareto front. We will also further study if this correlation can be *user-dependent* or *group-dependent*.

We start by computing what we call a *similarity heat map* of the user ratings normalized across all subjects and across all features combinations. The goal is to produce a generic and dithered view of the optimization space in order to have a global view of which portions of this space are rated higher than others. In order to obtain this overview, we developed a specific procedure that creates a 10×10 *similarity grid* matrix which partitions the space into equally-sized bins in both dimensions. We aggregate solutions from every combination in the corresponding bins depending on their coordinates in the optimization space. This aggregation is performed based on a centroid rule between the solutions and the closest points in the grid. Finally, we take the means of ratings for all combinations, users and solutions in their respective bins in the space (each bin is normalized separately depending on its cardinality). The results of this analysis are presented in Figure 45.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

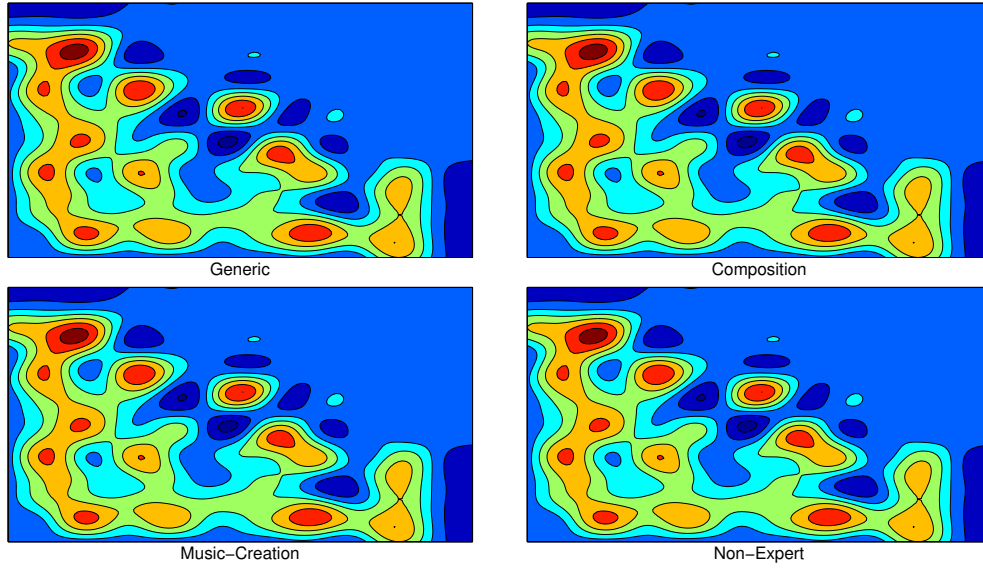


Figure 45: *Similarity heat maps* representing the weight of similarity ratings for all subjects and group-wise subjects independently of the features involved.

[[SPATIAL HOTSPOTS ANALYSIS]]
 [[COMPARISON MONO / MULTI]]
 [[VARIANCE-BASED ANALYSIS]]
 [[GROUP-WISE ANALYSIS]]

Based on the same analysis procedure, we also use the *similarity grid* to provide the kernel density estimates, as presented in Figure 46.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]
 [[SPATIAL HOTSPOTS ANALYSIS]]
 [[COMPARISON MONO / MULTI]]
 [[VARIANCE-BASED ANALYSIS]]
 [[GROUP-WISE ANALYSIS]]

To further deepen our analysis, we perform a *Mantel spatial correlogram* of the similarity ratings depending on their relative positions in the multiobjective space. The Mantel spatial correlogram tries to estimate the spatial dependences between distance classes. Therefore, this analysis allows to exhibit if the similarity ratings are correlated depending on their spatial relationships. Furthermore, the correlogram exhibits which classes of distances are significant and therefore we can also analyze if proximity in the solutions implies correlated ratings, ie. if the multiobjective similarity function can be explained by smooth transitions between nearby solutions.

The results of this analysis are presented in Figure 47. The correlogram is represented by a graph in which spatial correlation values (Mantel statistics) are plotted, on the ordinate, as a function of the geographic distance between solutions along the abscissa. The distance classes are discrete groupings computed to divide the elements in sixteen groups. The region-wide half-similarity between sites forms the reference line (zero abscissa). Positive Mantel statistics correspond to positive spatial autocorrelation, ie. distance classes which exhibit a significant correlation. On the other hand, negative Mantel statistics reveals distance classes at which ratings correlation is no more similar than random chance across the region. We test the Mantel statistic through a permu-

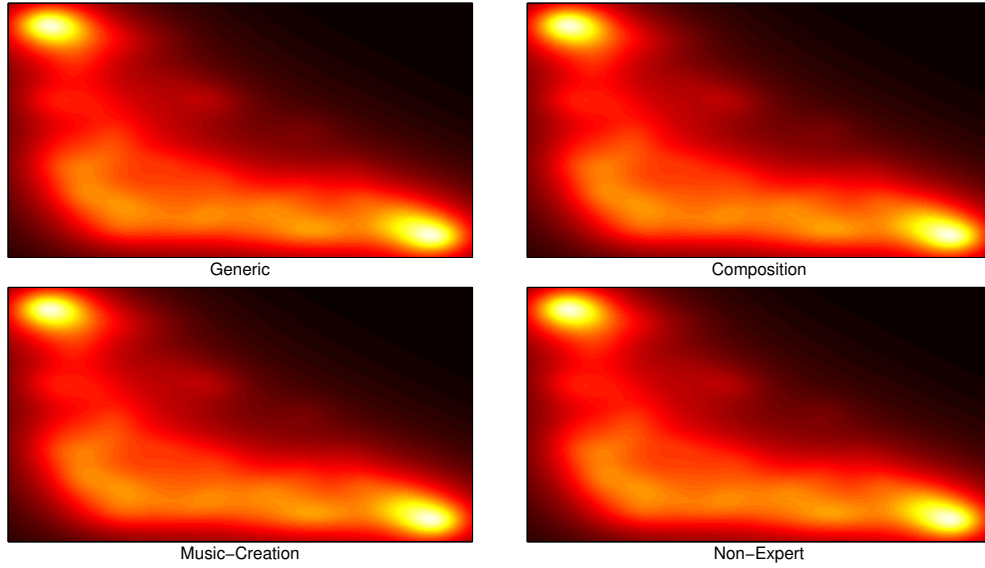


Figure 46: Kernel density estimates of similarity ratings distribution over the optimization space. The estimates takes ratings for all subjects and group-wise subjects independently of the features involved.

tational test and correction for multiple testing is applied by using Sturge’s rule and Holmes correction.

[[ANALYSIS OF OVERALL DISTRIBUTION]]
 [[SIGNIFICANT ANALYSIS]]
 [[NON-SIGNIFICANT ANALYSIS]]
 [[SPATIAL “SMOOTHNESS” ANALYSIS]]
 [[GROUP-WISE ANALYSIS]]
 [[CONCLUSIONS ON DIRECTIONS OF LISTENING]]

8.4.3 Usability evaluation of audio querying paradigms

As stated previously, we divide the evaluation of usability for our novel audio querying paradigms between their three main components, namely *efficiency*, *effectiveness* and *satisfaction*. In these analyses, we combine *quantitative* and *qualitative* measures to both know *which* behavior occurred and have some idea *why* it occurred.

Efficiency

The [Tsk-2] task requires for users to rate the generic sound similarity of results based purely on the sound information. It gives an opportunity to evaluate the overall *efficiency* of the search algorithm by using the similarity ratings in a straightforward fashion. Therefore, we first evaluate the *overall pertinency scores* averaged over every participants. We analyze these ratings for all combinations depending on the group and the features used. The results are presented in Table 8. This table presents the results for average and group-wise similarity ratings depending on the features involved.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]
 [[COMPARISON OF FEATURE SCORES]]

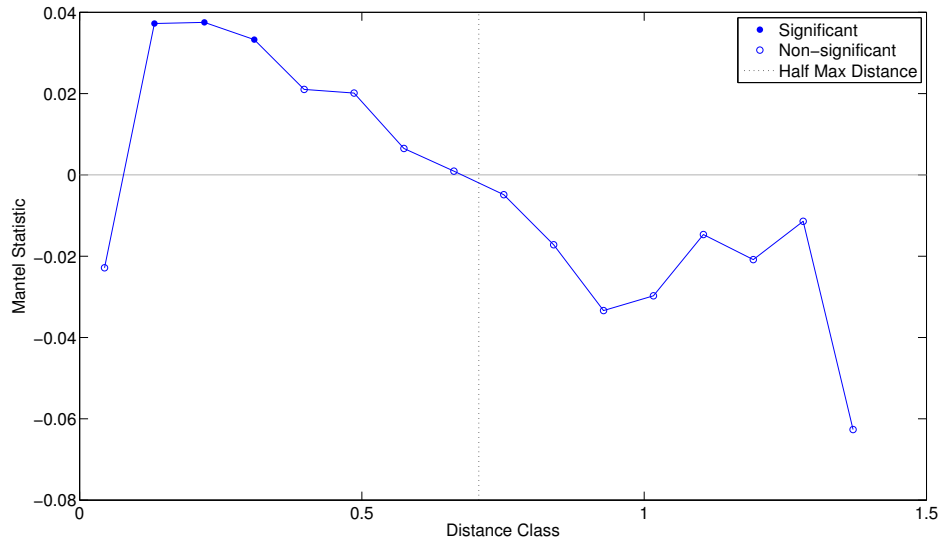


Figure 47: Mantel spatial correlogram performed over the similarity ratings depending on their distances to each other.

	Overall		Composers		Producers		Non-Experts	
	μ	σ	μ	σ	μ	σ	μ	σ
Average	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Fundamental	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Inharmonicity	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Loudness	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Noisiness	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Centroid	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Skewness	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2
Spread	5.6	1.2	5.1	0.2	5.1	0.1	5.1	0.2

Table 8: Overall similarity ratings for the pertinency task depending on the feature and the users groups.

[[VARIANCE-BASED ANALYSIS]]**[[FEATURE-DEPENDENT ANALYSIS]]**

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]*Effectiveness*

The last tasks were designed to assess the *effectiveness* (and *satisfaction*) of the MOSEQ and QVI paradigms, ie. to show that these audio querying paradigms can allow for efficient and intuitive retrieval. Therefore, we want to see if these systems allow users to easily find specific sounds with precise pre-defined structures. The fact that these sounds are imposed allows to normalize the results on effectiveness between various users (by removing the variability induced by *free retrieval* situations). We analyze the success of users in finding sounds, as well as their behavior during this process. In order to study the effectiveness, we use several performance measures. First, we use the *task completion time* (*mean*, *cumulative* and *query-dependent*) which give us an idea of how fast users can find a specific sound. To analyze the user behavior, we process the *number of queries* made by the user, *number of features modification*, *number of audio samples found* by the system against the *number of audio played* by the user and finally also count the *total amount of audio played*. We also provide a *performance histogram* on the *query specification* to see if there is a learning curve in user behavior, ie. does the user find solutions faster after getting used to the query system and thus improve their query formulation skills. This is done by analyzing the *number*, *type* and *duration of operations* but also the number of *coherent* (+) and *incoherent* (-) *query modifications* performed by subjects (ie. does the query modifications improve or worsen the results).

MOSEQ We start by evaluating the effectiveness of the MOSEQ paradigm, for which the overall results are summarized in Table 9.

[[TASK COMPLETION TIME ANALYSIS]]**[[NUMBER OF QUERIES ANALYSIS]]****[[NUMBER OF AUDIO PLAYED ANALYSIS]]****[[NUMBER OF FEATURE MODIFICATIONS ANALYSIS]]****[[FEATURE-DEPENDENT ANALYSIS]]**

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]**[[OVERALL EFFICIENCY AND USABILITY]]****[[TASK COMPLETION TIME ANALYSIS]]****[[NUMBER OF QUERIES ANALYSIS]]****[[NUMBER OF AUDIO PLAYED ANALYSIS]]****[[NUMBER OF FEATURE MODIFICATIONS ANALYSIS]]****[[CONCLUSION ON MOSEQ EFFECTIVENESS]]**

QVI We now evaluate the effectiveness of the QVI paradigm. Results are presented in Table 10

[[TASK COMPLETION TIME ANALYSIS]]**[[NUMBER OF QUERIES ANALYSIS]]****[[NUMBER OF AUDIO PLAYED ANALYSIS]]****[[NUMBER OF FEATURE MODIFICATIONS ANALYSIS]]**

	List approx	MOSEQ							
		Generic		Composers		Producers		Non-Experts	
		μ	σ	μ	σ	μ	σ	μ	σ
Task completion									
Mean	??	40.61	20.24	40.61	20.24	40.61	20.24	40.61	20.24
Cumulative	??	852.7	121.4	852.7	121.4	852.7	121.4	852.7	121.4
Queries required									
Mean	-	2.38	0.75	2.38	0.75	2.38	0.75	2.38	0.75
Cumulative	-	37.67	11.42	37.67	11.42	37.67	11.42	37.67	11.42
Audio played									
Number		2.38	0.75	2.38	0.75	2.38	0.75	2.38	0.75
Time		37.67	11.42	37.67	11.42	37.67	11.42	37.67	11.42
Features modification									
Mean	-								
Cumulative	-								
Feature-dependent									
Fundamental									
Inharmonicity									
Loudness									
Noisiness									
Centroid									
Skewness									
Spread									

Table 9: Results of the effectiveness of the MOSEQ paradigm evaluated through the *mean* and *cumulative* statistics for the *task completion* time, *number of queries required* and *number of features modifications* required as well as the *mean number of audio files played* and the corresponding *time required* to play these files. The last part of this table displays the *feature-dependent mean task completion time*.

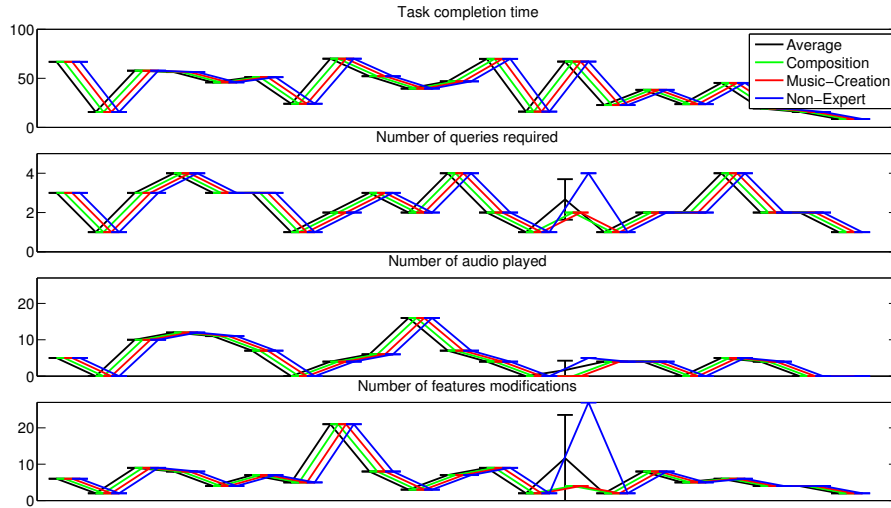


Figure 48: Query-dependent analysis of the main *effectiveness* measures for the MOSEQ constrained retrieval task

	List approx	QVI							
		Average		Composers		Producers		Non-Experts	
		μ	σ	μ	σ	μ	σ	μ	σ
Task completion									
Mean	??	40.6	20.24	40.61	20.24	40.61	20.24	40.61	20.24
Cumulative	??	852.7	121.4	852.7	121.4	852.7	121.4	852.7	121.4
Queries required									
Mean	-	2.38	0.75	2.38	0.75	2.38	0.75	2.38	0.75
Cumulative	-	37.6	11.42	37.67	11.42	37.67	11.42	37.67	11.42
Audio played									
Number		2.38	0.75	2.38	0.75	2.38	0.75	2.38	0.75
Time		37.6	11.42	37.67	11.42	37.67	11.42	37.67	11.42
Features modification									
Mean	-								
Cumulative	-								
Descriptor-dependent									
Fundamental									
Inharmonicity									
Loudness									
Noisiness									
Centroid									
Skewness									
Spread									

Table 10: Results of the effectiveness of the QVI paradigm, evaluated through the *mean* and *cumulative* statistics for the *task completion time*, *number of queries required* and *number of features modifications* required as well as the *mean number of audio files played* and the corresponding *time required* to play these files. The last part of this table displays the *descriptor-dependent mean task completion time*.

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - *Fundamental Frequency* - *Loudness* - *Noisiness* - *Spectral Centroid* - *Spectral Skewness* - *Spectral Spread*

[[GROUP-WISE ANALYSIS]]

[[OVERALL EFFICIENCY AND USABILITY]]

[[TASK COMPLETION TIME ANALYSIS]]

[[NUMBER OF QUERIES ANALYSIS]]

[[NUMBER OF AUDIO PLAYED ANALYSIS]]

[[NUMBER OF FEATURE MODIFICATIONS ANALYSIS]]

[[CONCLUSION ON QVI EFFECTIVENESS]]

Satisfaction

Analysis of the data collected through user surveys allows to perform a quantitative assessment of the *user satisfaction* towards the different components of the system. These surveys focus on the overall impression of subjects on the proposed paradigms. Users were asked to use a five-point Likert scale in order to rate their perceived *task difficulty*, *system utility*, *system satisfaction*, *paradigm usefulness* and potential *interest for re-using* the querying paradigms. Finally, the *perceived feature utility* ratings allow to divide the impression of subjects towards the *representation* of results in a multi-dimensional space, *quality* of the search results and *utility* of using the *temporal information* for matching.

MOSEQ We present the results of user satisfaction surveys for the MOSEQ paradigm in Figure 11.

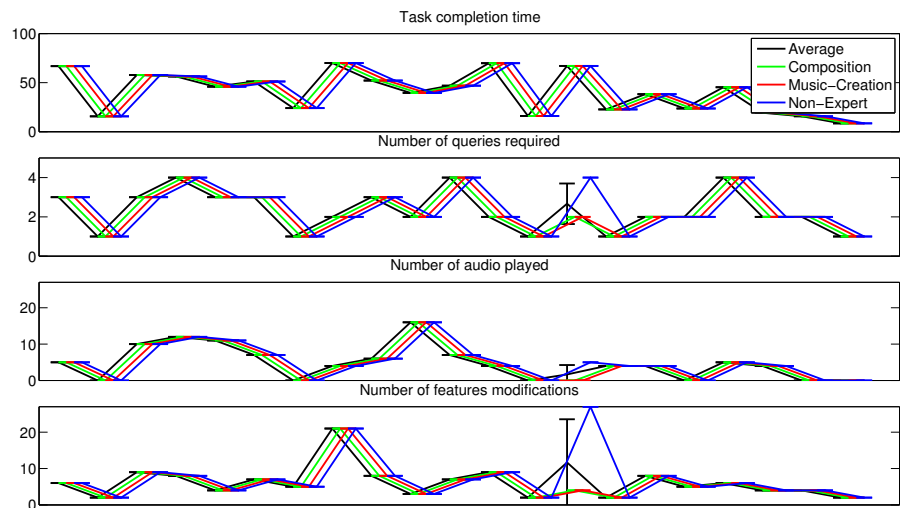


Figure 49: Query-dependent analysis of the main *effectiveness* measures for the MOSEQ constrained retrieval task

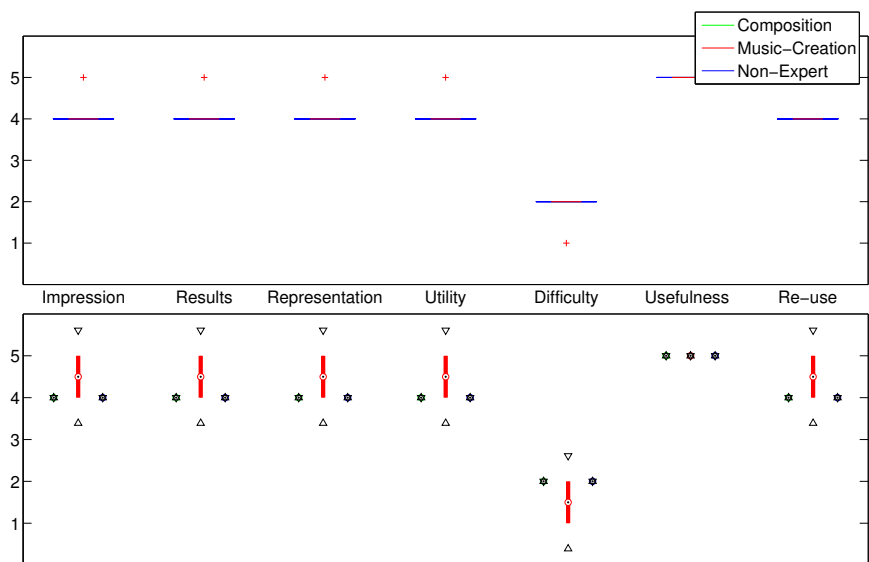


Table 11: Results of user satisfaction survey for the MOSEQ constrained retrieval task

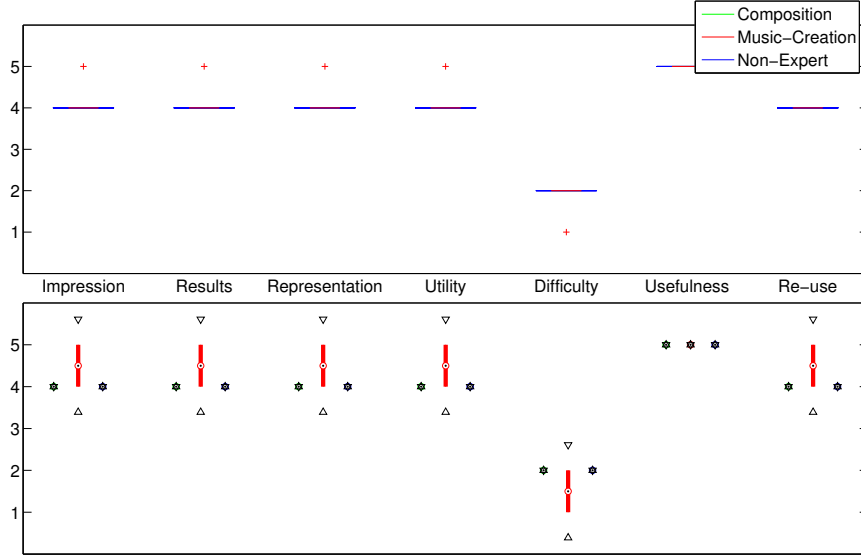


Table 12: Results of user satisfaction survey for the QVI constrained retrieval task.

QVI We present the results of user satisfaction surveys for the QVI paradigm in Figure 11.

Comparing the MOSEQ and QVI paradigms

After studying both paradigms separately, we try to compare the effectiveness and satisfaction between the MOSEQ and QVI paradigms in audio retrieval tasks. As our previous comparisons were based on theoretic points deriving from the usage of a list-based interface, we keep this point as reference and try to oppose the paradigms to each other in the same fashion.

EFFECTIVENESS [[COMPARISON Between Table 9 and Table 10]]

[[TASK COMPLETION TIME ANALYSIS]]

[[NUMBER OF QUERIES ANALYSIS]]

[[NUMBER OF AUDIO PLAYED ANALYSIS]]

[[NUMBER OF FEATURE MODIFICATIONS ANALYSIS]]

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[GROUP-WISE ANALYSIS]]

[[OVERALL EFFICIENCY AND USABILITY]]

SATISFACTION [[COMPARISON BETWEEN SURVEYS]]

8.4.4 Vocal control of spectral features

The QVI task offers a unique opportunity to analyze the extent of control over the human voice that subjects can exhibit for each particular feature. We consider that inference from the drawn queries gives us a valuable information on the vocal control.

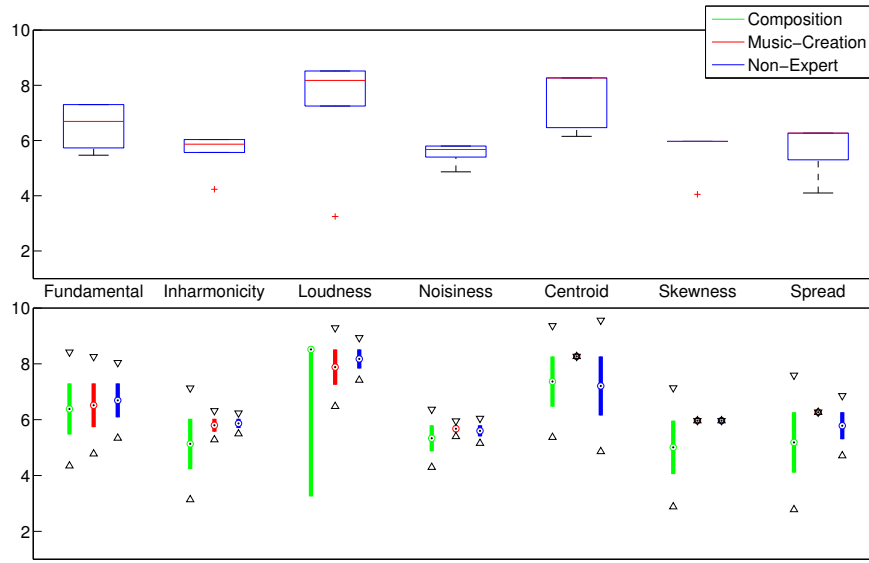


Figure 50: Distributions of mean vocal control strength for each sound feature, independently of the combinations used. Control scores for all subjects (up) and group-wise control (down).

As proposed for the directions of listening in the MOSEQ constrained retrieval task, we can see how similar the vocal imitations are to the actual feature of the original sound. However, like the previous task, we need to find a way to extract these directions of listening. Once again, we considering the time series distance between imitation and original features as a measure of *vocal control strength*. We therefore study the distances between the vocal shapes input by users and the true sound features computed over the queries, by using the *Dynamic Time Warping (DTW)* distance measure. We then normalize these distances and compute the mean vocal control strength per subject for each feature. The distribution of subjects and group-wise control strengths are presented in Figure 50.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

[[COMPARISON OF FEATURES CONTROL]]

[[VARIANCE-BASED ANALYSIS]]

[[GROUP-WISE ANALYSIS]]

As two features are required for each query, we can also apply these measures in a tournament fashion (as performed previously). We apply the same computation steps and presents the results in Table 13.

[[ANALYSIS OF OVERALL DISTRIBUTIONS]]

[[SCORE-BASED ANALYSIS]]

[[EIGENVECTOR-BASED ANALYSIS]]

[[VARIANCE-BASED ANALYSIS]]

[[FEATURE-DEPENDENT ANALYSIS]]

Inharmonicity - Fundamental Frequency - Loudness - Noisiness - Spectral Centroid - Spectral Skewness - Spectral Spread

	Fund.		Inha.		Loud.		Nois.		Centro.		Skewn.		Spread	
	Val	#	Val	#	Val	#	Val	#	Val	#	Val	#	Val	#
Tournament	6.53	3	5.60	5	6.53	3	5.60	5	6.53	3	5.60	5	5.60	5
Composers	6.38	3	5.13	5	6.38	3	5.13	5	6.38	3	5.13	5	5.13	5
Music-Creation	6.51	3	5.80	5	6.51	3	5.80	5	6.51	3	5.80	5	5.80	5
Non-Experts	6.69	3	5.87	5	6.69	3	5.87	5	6.69	3	5.87	5	5.87	5
Score-based	5.45	5	-7.16	3	5.45	5	-7.16	3	5.45	5	-7.16	3	-7.16	3
Composers	7.75	5	-8.4	1	7.75	5	-8.4	1	7.75	5	-8.4	1	-8.4	1
Music-Creation	1.65	5	-8.1	3	1.65	5	-8.1	3	1.65	5	-8.1	3	-8.1	3
Non-Experts	6.95	3	-5.0	5	6.95	3	-5.0	5	6.95	3	-5.0	5	-5.0	5
Eigenvector	0.0	1	-0.0	6	0.0	1	-0.0	6	0.0	1	-0.0	6	-0.0	6
Composers	0.32	1	0.32	3	0.32	1	0.32	3	0.32	1	0.32	3	0.32	3
Music-Creation	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	1	-0.75	3	-0.75	3
Non-Experts	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2	0.2	2
Variance	0.94	1	1.56	2	0.94	1	1.56	2	0.94	1	1.56	2	1.56	2
Composers	0.94	1	1.88	2	0.94	1	1.88	2	0.94	1	1.88	2	1.88	2
Music-Creation	0.94	1	1.58	2	0.94	1	1.58	2	0.94	1	1.58	2	1.58	2
Non-Experts	0.94	1	1.52	2	0.94	1	1.52	2	0.94	1	1.52	2	1.52	2

Table 13: Tournament-based analysis of listening directions for the constrained retrieval task.

[[GROUP-WISE ANALYSIS]]

8.4.5 *Impact of skills*

Given the observed distributions of similarity ratings, we try to find potential correlations between the different skills (especially in feature knowledge) and the ratings provided by corresponding users. Each rating being feature-dependent, we have to study a bi-dimensional matrix $[S \times F]$ of mean ratings with S the number of subjects and F the number of features. We therefore perform a Canonical Correlation Analysis (CCA) on this dataset. The results are presented in Figure 51.

[[VARIANCE EXPLAINED - OVERALL ANALYSIS]]

[[ANALYZE FIRST CANONICAL AXIS]]

[[ANALYZE SECOND CANONICAL AXIS]]

[[SKILL-WISE ANALYSIS]]

[[USER-WISE ANALYSIS]]

8.5 GENERALIZATION

We now attempt to generalize our results and provide a higher-level interpretation of the computed analyses by trying to learn from the evaluation and most of all try to exhibit the weakness of our proposals. We can use these weaknesses to drive our evaluation mostly by learning from the qualitative data but also from the quantitative data on certain parts of the interaction process.

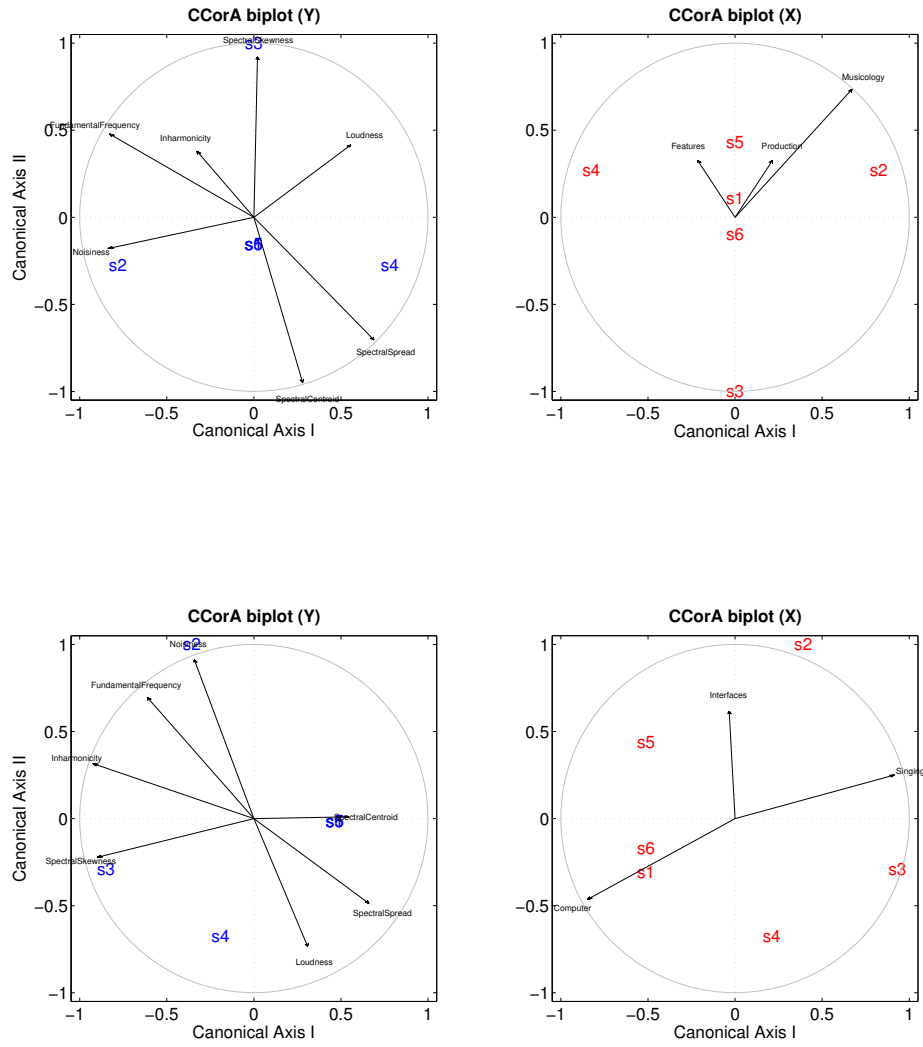


Figure 51: Canonical Correlation Analysis (CCA) between the user-wise similarity tournament matrix and the user self-rated skills.

9

CONCLUSIONS OF THIS PART

Motivated by observations in the field of audio similarity matching, we introduce the generic *MultiObjective Time Series* (MOTS) matching problem. This problem can be applied to any problem in which several time series should be matched jointly, without favoring any dimension in the process. The goal of MOTS is therefore to provide a more flexible assessment of multiple time series similarities.

Equipped with the basic notions of time series matching and multiobjective optimization, we start by introducing the generic MOTS matching problem (Section 7.1), we further exhibit the core differences between this novel problem and multivariate matching (Section 7.2) and also quickly discuss its complexity. We then introduce two algorithms to solve this problem (Section 7.3) and show their efficiency on massive sets of data (Section 7.4). We further compare the efficiency of these algorithms between real and synthetic sets of data (Section 7.4.2). Finally, we discuss the application of the MOTS framework to the audio retrieval problematic. We show that thanks to this framework, we can easily construct two innovative audio querying paradigms (Section 7.5), that allow to go beyond the traditionnal audio query applications.

Part IV

HYPERVOLUME CLASSIFICATION (HV-MOTS)

10

HV-MOTS CLASSIFICATION

As we have just seen in the previous chapter, the MOTS approach allows to find the set of efficient solutions in a database given multiple time series features. We showed that this approach gives access to more flexible sets of solutions in the context of retrieval and querying. However, we could wonder if this flexibility can also benefit to other field of studies. Notably, as the MOTS framework is intended to find efficient sets of similarities in a database, it fits the basic requirements of classification paradigms. However, given the definition of Pareto optimality, there is normally no way of ranking the different elements between each other. Therefore, there is no straightforward criterion to make a final classification decision. We introduce the notion of *hypervolume dominated* by a class and show how to use it as a classification criterion. Our main idea is that by not merging every dimensions into a single distance measure, we can benefit from a more accurate and flexible view on the properties of various classes. We will show that using this multi-objective flexibility with the adequate hypervolume criterion allows to construct a classifier that significantly outperforms state-of-art methods. Specifically, our approach exhibits a statistical superiority over the 1NN-DTW classifier which has been consistently shown to still be the best performing classification scheme for time series Gudmundsson et al. [155, 156], Islam et al. [187], Radovanovic et al. [315], Radovanović et al. [314].

10.1 MULTI-OBJECTIVE CLASSIFICATION

Based on the MOTS selection, we still need a criterion to make the final classification decision, ie. to select which class is the best match to a given input. We introduce in this section three class selection criteria that allow to use the Pareto optimality and apply it to any classification problem. We consider that elements to be classified are compared to known items thanks to different distance measures on several dimensions. We also consider that these distances are not merged and give access to a complete multidimensional distance matrix.

Pareto cardinality

Given the Pareto set, we can first simply look at its cardinality. Therefore, our first selection criterion can be obtained by counting the number of occurrences of each class c in the Pareto front \mathcal{P} . The selected class \mathcal{C}_s is the most represented in the front.

$$\mathcal{C}_s = \underset{c}{\operatorname{argmax}} (|\{p_i \in \mathcal{P}, \text{class}(p_i) = c\}|) \quad (10.1)$$

This is obviously a simple criterion and we can expect it to be less efficient in higher dimensions. We term this method *MOTS* in the following.

Nearest Pareto

Given an input, each class c can provide a different Pareto front \mathcal{P}_c depending on its distances sub-matrix. Therefore, our second selection criterion can be obtained by

computing the Pareto front of each class and then computing the mean distance between the input and all elements in the front. The selected class is therefore the one which implies the minimal distance.

$$\mathcal{C}_s = \underset{c}{\operatorname{argmin}} \left(\frac{1}{n} \sqrt{\sum_{i=1}^n \|p_i\|^2}, p_i \in \mathcal{P}_c \right) \quad (10.2)$$

We term this method *Nearest Pareto MOTS (NP-MOTS)* in the following.

Hypervolume domination

We now introduce a novel criterion based on *hypervolume* domination. This measure has been used in multi-objective optimization with Genetic Algorithms (GA) [436] as a performance indicator, i.e. only to differentiate the quality of different algorithms. However, to our best knowledge, it has never been used as a classification criterion. The idea behind this measure is that every point in a multi-dimensional space, defines a hypervolume which indicates the portion of space dominated by this point. For a n -dimensional space, $n \in \mathbb{N}$, the hypervolume of a box in \mathbb{R}^n generated by two points $a = (a_1, \dots, a_n)$ and $b = (b_1, \dots, b_n)$ is defined as

$$\mathcal{H}(B) = \prod_{i=1}^n (b_i - a_i) \quad (10.3)$$

The *hypervolume dominated* by a Pareto front \mathcal{P} given a reference point $r_p = (r_p^1, \dots, r_p^n)$ is given by the union of hypervolumes dominated by each point in the front

$$\mathcal{H}(\mathcal{P}) = \mathcal{H}\left(\bigcup_i B_i\right) = \mathcal{H}\left(\bigcup_{(p_1, \dots, p_k) \in \mathcal{P}} [p_1, r_p] \times \dots \times [p_k, r_p]\right) \quad (10.4)$$

These notions are shown in Figure 52 (left). Point p_1 defines a box B_1 (darker gray) with the reference point r_p . Each point of this set also implies a corresponding domination box. The hypervolume dominated by the Pareto front is therefore the union of hypervolumes dominated by each point in the front. Figure 52 (right) shows the benefits of this measure when comparing two distributions. Even though the first class have more elements belong to the final Pareto front, its dominated hypervolume \mathcal{H}_1 is smaller than the hypervolume \mathcal{H}_2 of the second class. Therefore, the hypervolume indicates both the fitness of a distribution and its spread over the optimization space. Furthermore, compared to a NN or NC rule, it summarizes the behavior of the whole class with respect to the input rather than the position of the input relative to a single element of the classes. In our implementation, we use the hypervolume computation algorithm proposed by [130]. We compute the hypervolume dominated by the Pareto front of each class. The selected class is therefore the one which induces the largest dominated hypervolume.

$$\mathcal{C}_s = \underset{c}{\operatorname{argmax}} (\mathcal{H}(\mathcal{P}_c)) \quad (10.5)$$

We term this approach *HyperVolume-MOTS (HV-MOTS)* in the following.

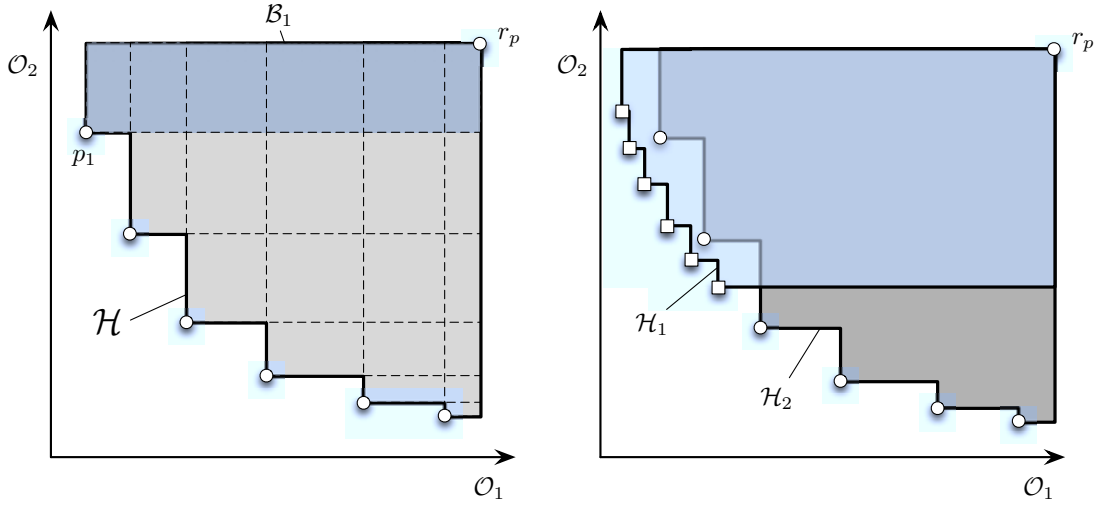


Figure 52: (Left) Hypervolume dominated by a Pareto front given the reference point r_p in a 2-dimensional space. The darker gray subpart defines the box B_1 which is dominated by point p_1 . The hypervolume \mathcal{H} dominated by the Pareto front is defined as the union of all boxes dominated by each point of the front. (Right) Comparison of two dominated hypervolumes \mathcal{H}_1 and \mathcal{H}_2 . Even though the first class have more elements belong to the final Pareto set, its hypervolume \mathcal{H}_1 is smaller than the hypervolume \mathcal{H}_2 of the second class.

10.2 DISTANCE-BASED CLASSIFIERS

The HV-MOTS classification framework falls in the category of distance-based classifiers. To see the novelty implied by our proposal, we start by comparing it to other distance-based classifiers. Figure 53 illustrates these concepts. The element to be classified is represented by the cross at the origin of the space. The *Nearest-Neighbors* techniques will try to find the nearest element(s) based on the norm of the distance vector, thus defining the selected class accordingly. The *Nearest Center* technique performs the same analysis but by first computing the centroid of each class and then selecting the nearest one. The *MOTS* paradigm computes the Pareto front and then selects the class that is the most represented inside the set. The *NP-MOTS* paradigm first computes the Pareto set of each class and then selects the nearest one based on their mean distances. Finally, the *HyperVolume-MOTS* technique computes the hypervolume dominated by each class and then selects the class with the largest hypervolume.

As we can see in this figure, unlike the other classification schemes, the HV-MOTS technique does not perform a linear merging of the distance measures in every objectives. Instead, the computation of the hypervolume allows to account for two aspects of the class distributions. First, the *proximity* (as for the other schemes) is implied, given that the hypervolume will grow if the elements of the front are closer to the input. However, this criterion also accounts for the *spread* of the classes over the distance space. This therefore allows for flexible matching over both dimensions separately. Hence, it analyzes if classes to perform well in *every* directions of optimization.

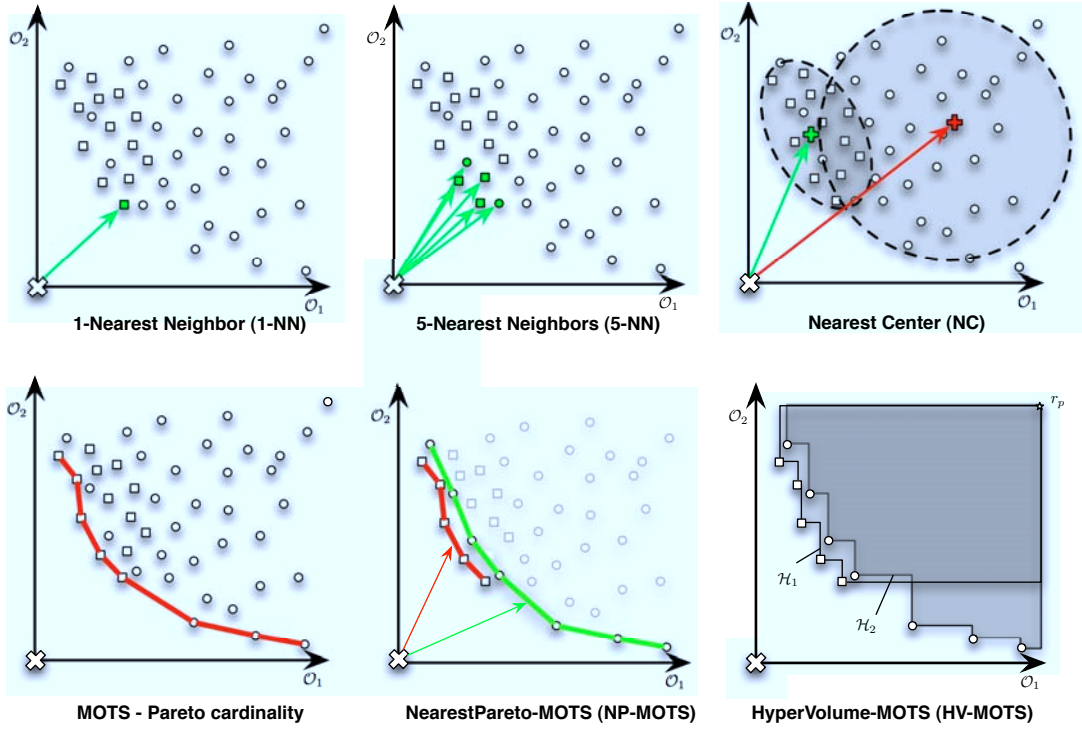


Figure 53: Comparison of distance-based classifiers. The element to be classified is represented by the cross at the origin of the space. The *Nearest-Neighbors* techniques select the class of nearest elements based on the norm of their distance vector. The *Nearest Center* technique first computes the centroid of each class and then selects the nearest one. The *MOTS* paradigm computes the Pareto front and then selects the most represented class. Finally the *HV-MOTS* technique computes the *hypervolume dominated* by each class and then select the class with the largest one.

10.3 COMPARISON TO OTHER CLASSIFIERS

We provide here a brief comparison on the class boundaries provided by the HV-MOTS classification framework, as opposed to other well-studied state-of-art classifiers. We separate this comparison between the methods which provide *linear* or *non-linear* class boundaries. We underline that this comparison is far from exhaustive but we selected a subset of classification methods that will be used in the subsequent large scale analysis of the HV-MOTS classifier. Figure 54 illustrates the comparison of several classifiers including our approach in terms of class boundaries in *feature* space.

NEAREST-CENTER As discussed previously, the NC classifier works by finding the center of class distributions and, therefore, the most likely positions of their elements. Then the input is compared only to these centroids. As we can see in Figure 54, the corresponding boundary is *linear* as the distance function delineates the boundaries between class pairs.

NEAREST-NEIGHBOR The *1-NN* approach tries to find the element that is closest to the input in terms of features distances, usually through their Euclidean norm. Hence, as it compares an input to every element of a class, the reference point changes

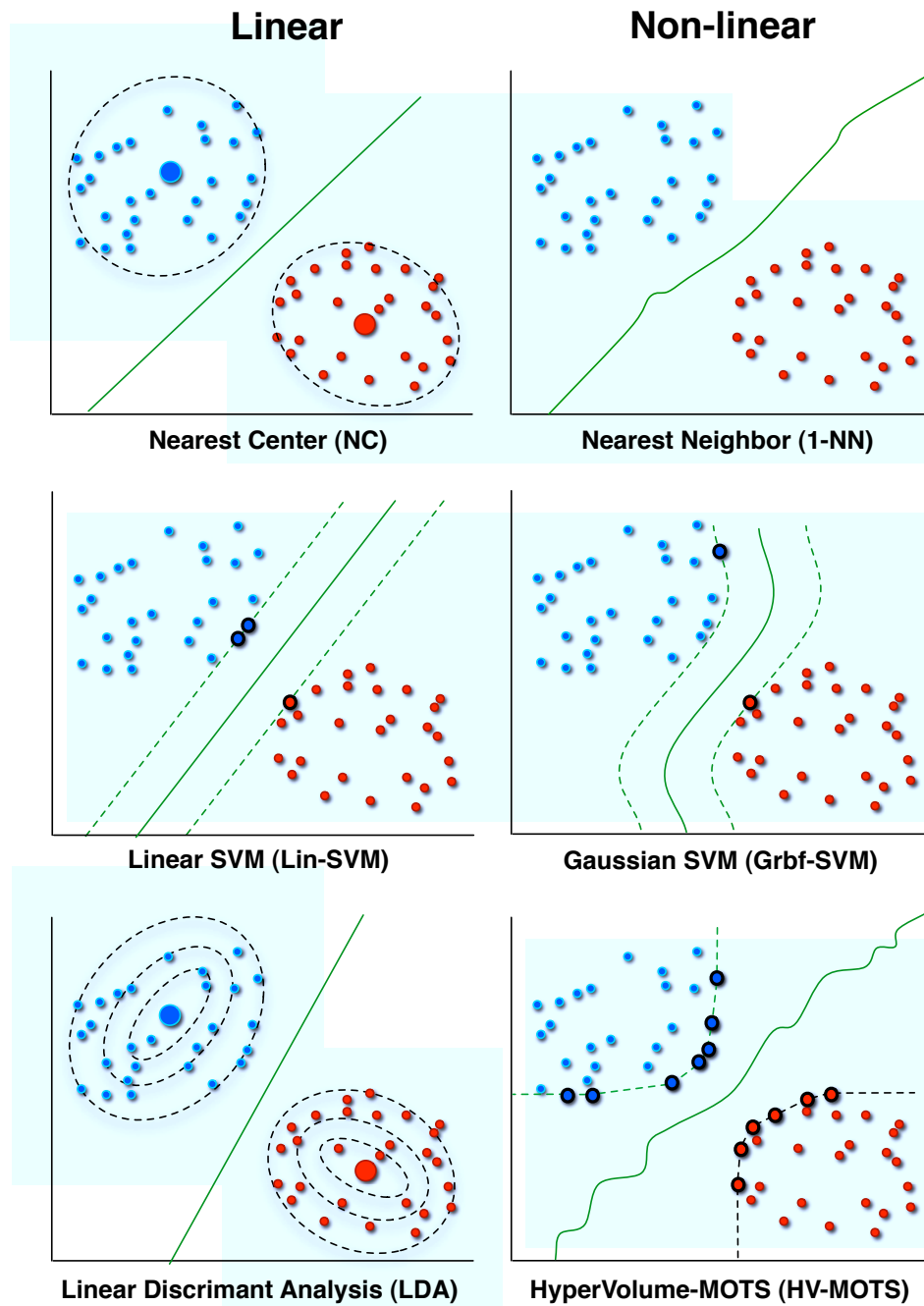


Figure 54: Comparison of several classification approaches based on the class boundaries that they define in *feature* space. The techniques are separated between their definition of *linear* (left) or *non-linear* (right) class boundaries.

depending on the input. Therefore, as we can see in Figure 54, the resulting boundary is *non-linear*.

SUPPORT VECTOR MACHINES The SVM classification is based on two key ideas. First, it relies on the notion of *maximal margin* between the class boundary (solid line in Figure 54) and the closest elements to it. This gives a second boundary (dotted line in Figure 54), which defines the *margin* for decision. Hence, the elements selected for this margin are called *support vectors*. The second notion lies in the concept of *kernel*. This idea is introduced to solve the classification cases where the different classes are not linearly separable. Therefore, the kernel is applied to transform the data into a higher dimensionality where the separability is linear. We detail two possible kernels for SVM.

Linear kernel The simplest case of SVM is to use a linear kernel (ie. no dimensionality modifications are applied). This allows to position a decision hyperplane between the classes. The algorithm selects the hyperplane that provides the maximal margin (maximum separation) between the classes. As we can see in Figure 54, the Lin-SVM provides a *linear* separation between classes.

Gaussian kernel If the data requires a non-linear separation, a kernel function can define an alternate high-dimensional space. The main idea here is to use a set of functions called *Gaussian Radial Basis Functions* (GRBF), which are applied to the datas. The resulting space allows to use a linear separation based on the same hyperplane selection. Figure 54 exhibits the corresponding *non-linear* class boundary, maximal margin and support vectors.

LINEAR DISCRIMINANT ANALYSIS The *Linear Discriminant Analysis* (LDA) tries to find the dimension which maximize the inter-class variance and intra-class coherence. The ratio between these two values (sometimes called *inertia ratio*) is maximized in order to find the *linear discriminant* dimension. The data is then projected on this dimension in order to place a classification threshold between the classes means. This allows to obtain a hyperplane boundary between classes which is orthogonal to this dimension. As we can see in Figure 54, this leads to a *linear* class boundary.

HYPERVOLUME-MOTS We try to provide the main differences and resemblances between our proposal and the presented classifiers. First, the HV-MOTS approach works in a fashion similar to NN and NC by trying to select the most similar class directly from the distance matrix. However, our approach is fundamentally different from these two methods, as it does not use a “direct distance” approach between elements. Indeed, the HV-MOTS scheme studies the overall distribution of each class separately. Therefore, it allows to analyze both the *distance* and *spread* of various classes. Moreover, by not using an Euclidean norm for elements-to-input distance and not merging dimensions into a single measure, each of the dimensions are treated separately. HV-MOTS also falls into the category of non-linear class boundaries such as NN and Grbf-SVM. Compared to the NN approach, however, we can see in Figure 54 that the delimited boundary between classes is somehow more complex. This comes from the fact that it uses information of the density and distance on dimensions separately. However, it is interesting to note here that for a single dimension, HV-MOTS is strictly equivalent to 1-NN classification. Indeed, as only one dimension is involved, the largest hypervolume is necessarily induced by the nearest element.

10.4 DISCUSSION

10.4.1 Past results

In the context of time series classification, it has been repeatedly proven that the 1-NN classifier is extremely hard to overpower [405]. Several published classification studies confirms that 1-NN classification is still the best performing classification scheme for time series retrieval [155, 156, 187, 315, 314]. Some authors even point out that “*while there have been attempts to classify time series with decision trees, neural networks, Bayesian networks, support vector machines etc., the best published results (by a large margin) come from simple nearest neighbor methods*” [109]. Even the SVM which is one of the most powerful classification scheme available appears to be *at most* statistically equivalent to 1-NN but usually performs worst [155]. This comes from the fact that time series are inherently high-dimensional data. Therefore, it is very difficult to define statistical methods that could avoid this *curse of dimensionality*. Therefore, the best performing methods are those based on distance comparisons.

As it has been repeatedly shown (in time series classification) that the 1-NN classifier outperforms other classification approaches, we will focus on comparing our approach to this framework. We will show in the large scale study that our novel approach statistically outperforms the 1-NN selection scheme.

10.4.2 Advantages and drawbacks

We try to discuss here the theoretical advantages and drawbacks of the HV-MOTS classification framework. As this method is a distance-based classifier, it provides similar disadvantages. First, as the NN classifiers, an obvious disadvantage is the time complexity of making the final decision. Indeed, the distance-based methods requires to compute the complete distance matrix. For the HV-MOTS selection, we can use the algorithms introduced in Section 7.3 that provides a significant speedup to find the Pareto front of each class. However, this imply to construct a different index for each class.

A major advantage of the HV-MOTS method is that *no assumptions* are made on the data or its distribution. Therefore, it can classify any type of data. Furthermore, as other distance-based methods, *no training* is required by this approach. Therefore, there is no need for extensive statistical models to be constructed beforehand. Another advantage of distance-based methods is that they are not bound by the *curse of dimensionality* involved in time series matching. Therefore, through the use of pairwise distances, the underlying data can be of any complexity without requiring the application of dimensionality reduction techniques. Furthermore, the HV-MOTS approach requires *no parameters* for classification, which makes it applicable to any dataset in a straightforward manner.

Finally, it has been shown that one of the main problem of k -NN techniques is that it tends to be very sensitive to irrelevant features as every dimensions contribute to the classification. A way of avoiding this problem is to weigh the different dimensions [97]. Therefore, compared to 1-NN selection, this is the main enhancement provided the HV-MOTS approach. Indeed, it specifically targets this problem and allows to alleviate its drawbacks, by considering *every possible weighting* in 1-NN selection. This theoretical improvement is straightforward from the fact that each element in the Pareto front is the best 1-NN selection given a particular set of weights.

10.4.3 Comparison

We try to provide here a deeper comparison of the HV-MOTS classification and 1-NN selection based on the class boundaries induced by both methods in feature space. Figure 55 illustrates the comparison of boundaries between 1-NN and HV-MOTS on two synthetic sets of data. The first problem (up) represents synthetic data where the classes are *almost* linearly separable. However, a slight intersection of class properties appears at the boundary. As we can see, the 1-NN defines harsh boundaries around these areas, by only taking into account the proximity of *nearest* elements. Oppositely, the HV-MOTS approach also accounts for the distribution of other points in the surroundings. Therefore, the overall boundary defined by HV-MOTS seems to be smoother than the 1-NN boundary. This fact appears more clearly in the second problem (down) where the set of classes is completely mixed and, therefore, with no linear separation. We underlined some areas of relevance over the feature space. For instance, the bottom left area is of particular interest. As we can see, the blue point appears to be an outlier in a region which has an higher density of red points. The 1-NN selection will only define the boundaries based on proximity. Oppositely, the HV-MOTS method clearly exhibits this outlier nature as it accounts for the *spread* and, consequently, the density of local points.

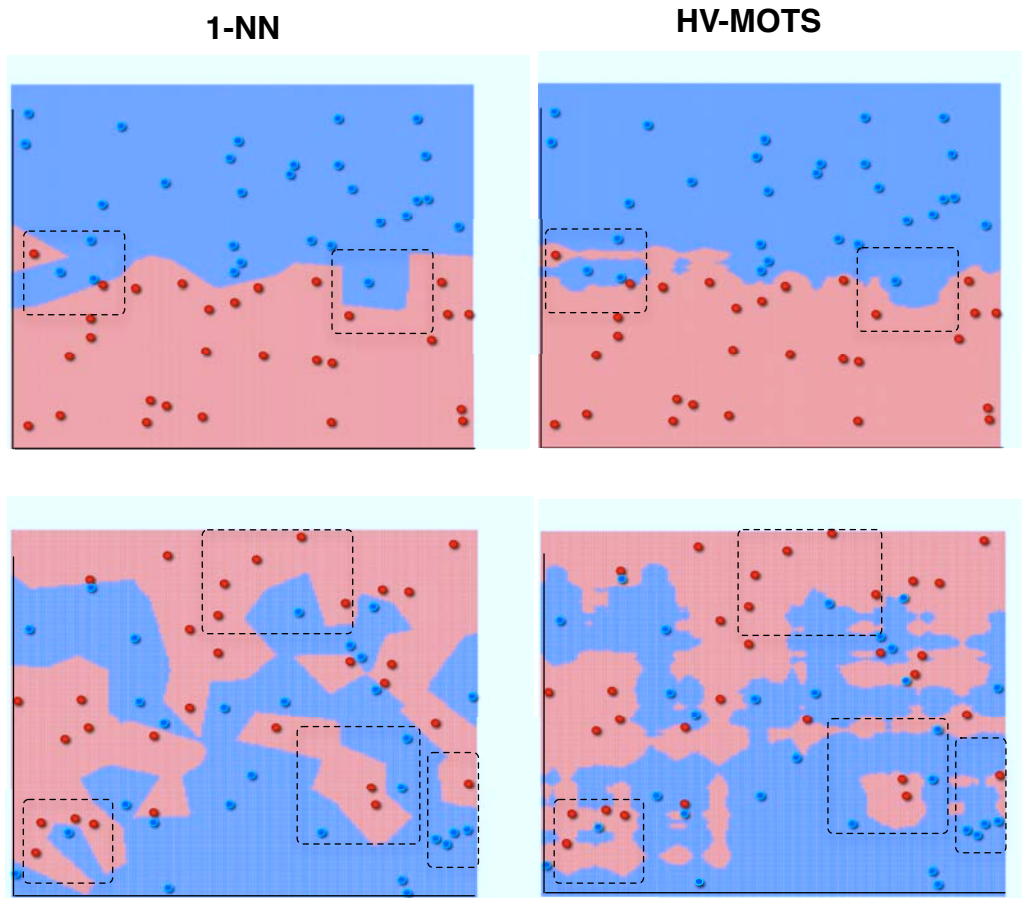


Figure 55: Comparison of the classification boundaries represented in feature space implied by 1-NN selection or HV-MOTS classification algorithms. The first problem (up) represents synthetic data where the classes are *almost* linearly separable. The second problem (down) represents a mixed set of classes data with no linear separation.

LARGE SCALE STUDY

In order to thoroughly evaluate the efficiency of the HV-MOTS classification framework, we follow several guidelines provided by past research in multiple classifiers comparison [105, 208, 271, 337, 338]. Therefore, we test our approach on a wide range of datasets that covers a variety of several scientific fields. This allows to extensively assess our proposal and at the same time provide conclusions that are free from any *data bias*. Moreover, we can also analyze in *which* situations the HV-MOTS classifier performs best. This entails the scientific field to which it is applied, but also the *characteristics* of the problem (number of classes, samples and objectives). Finally, this evaluation framework also offers an opportunity to analyze which features were chosen by the system to obtain the best classification accuracy. The analysis of this choice can provide some insights on the underlying classification problem.

Therefore, we collected several datasets that all meet the same requirements. First, they should be focused on *time series* data and therefore focus on processes that are inherently temporal. Secondly, they should be structured around a *classification* task. In addition, they should also have a *multidimensional nature* and, therefore, provide at least two objectives or features. Finally, these datasets should optimally be linked to research papers that provide state-of-art classification results. Therefore, we gathered corresponding accuracies, characteristics and methods used for further comparison.

11.1 DATASETS SUMMARY

We provide in this section a brief summary of the datasets collected for the evaluation. For the sake of clarity, we only give here an *overview* of the datasets. However, for the sake of completeness, the complete description of datasets along with technical informations is available in annex of this document (cf. Section A.1).

Overall, we collected 40 datasets which characterize a range of classification problems. These datasets come from fields of *speech recognition*, *handwriting analysis*, *character recognition*, *movement analysis*, *hand sign recognition*, *brain-computer interface*, *EEG analysis*, *MEG analysis*, *ECoG analysis*, *cardiology*, *medical surveillance*, *climatology*, *radar analysis* and *robotics*. The characteristics of these datasets exhibit wide differences. The number of features varies between 2 and 306 (mean : 40.75, median : 12), the number of classes varies between 2 and 183 (mean : 20.9, median : 8) and the number of samples varies between 80 and 164860 (mean : 7540.12, median : 1347).

Regarding the features themselves, their properties vary widely as well. This once again allows to abstract the potential conclusions from data bias. The corresponding time series span from 36 to 2134 time points which represent various underlying durations. The corresponding size of datasets range from 4.9 Mo to 1721 Mo (mean : 211.448 Mo, median : 83 Mo) for a complete size of 11572.04 Mo.

Datasets	Description	Features	Classes	Samples
Arabic digit [165]	Spoken arabic digits	13	10	8800
Artificial characters [53]	Character recognition	2	10	6000
Australian signs [199]	Hand sign recognition	10	95	6650
Australian signs (HQ) [200]	Hand signs (HQ)	20	95	2565
BciIII-01-Tubingen [313]	Brain-Computer	64	2	378
BciIII-02-Albany [318]	Brain-Computer	64	36	185
BciIII-03a-Graz [49]	Brain-Computer	60	4	840
BciIII-03b-Graz [237]	Brain-Computer	2	2	2760
BciIII-04a-Berlin [396]	Brain-Computer	118	2	1400
BciIV-01-Berlin [50]	Brain-Computer	64	3	1400
BciIV-03-Freiburg [50]	Brain-Computer	10	4	480
Biomag-2010 [382]	EEG analysis	274	2	780
Brain-Computer [101]	Brain-Computer	8	2	160
Challenge-2011 [407]	Cardiology	12	2	2000
Character-trajectories [297]	Character recognition	3	20	2858
Dachstein [158]	High altitude medicine	3	2	698
Digit-hands [9]	Character recognition	3	10	10992
Eeg-alcoholism [431]	Medical analysis	64	6	650
Eeg-epfl [179]	EEG analysis	34	36	26646
Forte [42]	Climatology	2	7	121
Gaitpdb [236]	Gait analysis	18	2	306
Handwritten [131]	Character recognition	2	183	8235
Ionosphere [118]	Radar analysis	34	2	358
Japanese-vowels [229]	Speech analysis	12	9	640
Libras [341]	Movement recognition	2	15	360
Meg mind reading [221]	MEG analysis	306	5	1330
Neuro-emotional [279]	ECoG analysis	128	8	300
Neuro-visual [279]	ECoG analysis	128	8	200
Pen Characters [307]	Character recognition	2	62	1364
Pen Characters [80]	Character recognition	2	97	11640
Person activity [203]	Movement analysis	12	11	164860
Physical action [376]	Movement analysis	8	20	80
Ptbdb [156]	Cardiology	15	9	2750
Robot failures [67]	Robotics	6	5	463
Slpdb [62]	Sleep apnea analysis	7	7	4085
Sonar [371]	Sonar analysis	60	2	208
Synemp [196]	Climatology	2	2	20000
Vfdb [228]	Cardiology	2	15	600
Vicon physical [376]	Physiological analysis	26	20	2000
Wall-robot [132]	Robotics	28	4	5460

11.2 METHODOLOGY

We intend to compare the classification results between the HV-MOTS framework and several classifiers. Therefore, we will outline in the next section the guidelines from previous research on the comparison of multiple classifiers and datasets. First, we underline the argument put forward by Keogh and Kasetty [209], that such comparisons should try to be free of both *implementation* and *data bias*.

Definition 32. *Implementation bias* is the conscious or unconscious disparity in the quality of implementation between a proposed approach and the competing approaches.

As proposed by [209], we tried to perform an extremely conscientious implementations of all approaches, combined with diligent explanations of the experimental process.

Definition 33. *Data bias* is the conscious or unconscious use of a particular set of testing data to confirm a desired finding.

In order to avoid data bias, we present the results from all the gathered datasets. These sets are from several scientific fields, with various properties, number of features, classes and samples in order to maximize the diversity of evaluation data.

11.2.1 Evaluation framework

In order to perform a comprehensive and thorough evaluation of the HV-MOTS classification scheme, we discuss here the indications provided by several researches focused on the comparison of classifiers over multiple datasets Demsar [105], Misaki et al. [271], Salzberg [337, 338]. That way, we try to avoid falling into the pitfalls of comparisons that could mislead the readers. As noted by Salzberg [337], very large databases are bound to contain some statistical anomalies. Hence, careful attention should be directed to these phenomena that could invalidate experimental comparisons. Therefore, we avoid the *multiplicity effect* and do not use the simple paired t-test, but rather use the non-parametric Friedman test, as also proposed by Demsar [105], with the corresponding post-hoc tests for comparing all classifiers over every datasets. Therefore, we first use Tukey-Kramer Honestly Significant Difference (HSD) test [96] over the results of Friedman's ANOVA, as proposed by Demsar [105] and also suggested in other fields of evaluation [114]. Finally, we present the *critical difference graphs* Demsar [105] which allows to exhibit the true statistical superiority and eventual groups of statistical equivalence between various methods. As pointed by Salzberg [337], the *repeated careful tuning* is usually performed in order to produce a desired result on a particular set of data. Hence, the results produced by this methodology might be misleading. Therefore, *we do not operate any form of tuning* nor do we make *any kind of pre-processing* on the datasets in order to enhance results. We underline here the fact that we do not perform any further extraction of features and work solely on the raw time series available. That way, we focus solely on the classification method used rather than the discriminative power of the underlying feature sets.

In order to compute the similarity between elements, we directly use the time series features in each dataset. We first normalize the series using *zero mean* and *unit variance* transformations. The mean and deviation of the time series are kept and used separately as features. Finally, to compare the time series, we use the Dynamic Time Warping (DTW) distance and the Euclidean distance computed on down-sampled series, which

leads to two different distance measures per time series. Therefore, we have four different features comparisons (mean, deviation, DTW and Euclidean) for each time series. This leads to a distance matrix computed between every elements of the dataset.

As our method does not require any training, we use the *Leave-One-Out* evaluation methodology. That is, each file is first withdrawn from the dataset and then input for classification with the remaining set acting as a database. In order to compare different methods, we perform large-scale experiments by testing combinatorial possibilities among every available feature for each dataset. Therefore, we start by performing classification with only one feature. We then test classification with every combination of two features, and so forth. Given that this testing methodology implies an exponentially growing number of tests, we keep only the top performing half of the features set after each step. The selection of retained features is based on their mean classification accuracies (across methods). We repeat this procedure and halve the set of available features (in which to choose the objectives for classification) until the number of remaining features is less than the current number of objectives. This testing methodology ensure *completeness* as we will test statistical significance independently of the underlying feature selection. Therefore, we test the classification power of each method on *every* subset of available features. Hence, we focus on the *classification criterion* instead of the set of features used. Furthermore, this methodology also allows to perform a *feature analysis*, by extracting for each dataset the best performing features combinations.

11.2.2 Hardware

As our testing methodology is computationnaly intensive, it required a high-performance supercomputing ressource in order to attain completion. We made all our calculations on the Guillimin cluster from the CLUMEQ supercomputing center of McGill university under the govern of Calcul Quebec and Calcul Canada. The exploitation of this supercomputer is financed by the Canadian fondation for innovation (FCI), the Research Concl of natural sciences (CRSNG), NanoQuebec, the RMGA and the Quebec research fund - Nature and technologies (FRQ-NT). The Guillimin server contains 1200 computing nodes and 34 infrastructure nodes. These nodes are all constituted by a pair of Intel Westmere-EP (Xeon X5650) processors each of which containing 6 processing cores and 24, 36 ou 72 gigaoctets of RAM memory. All nodes are linked together by a high performace Infiniband QDR network.

11.2.3 Algorithms implementation

In order to perform an exhaustive evaluation, we tried to implement several classification techniques. However, it should be noted that we are working with *time series* data which implies a different set of constraints than the usual classification problems. In time series classification, the best published results are usually provided by simple nearest neighbor methods [405]. As discussed earlier, the 1-NN method has been proven as the most efficient classifier over a wide variety of datasets [209]. Nevertheless, we also implemented the 5-NN, NC and SVM classifiers to ensure a thorough comparison. The 1-NN, 5-NN and NC classification techniques can be directly applied to the distance matrix. However, the SVM classifier usually requires a matrix of features for classification rather than a matrix of distances. Furthermore, we must face the high dimensionality of the datasets and can not use the time series features as direct input

to an SVM. Therefore, we follow the proposal of Gudmundsson et al. [155] and use a *proximity function* kernel which is designed to classify pairwise data.

- 1-NN We find the nearest element to the input (based on the norm of the distances) by computing the complete distance matrix and selecting the class accordingly.
- 5-NN The same idea applies to finding the five nearest element of the dataset and then selecting the corresponding class.
- NC The centroid of each class is computed based on the distance matrix. We then select the class with the nearest centroid.
- SVM As proposed by Gudmundsson et al. [155] we implement a SVM based on the *pairwise-proximity function* kernel (ppfSVM), which allows to directly use the distance matrix in order to perform classification.

11.2.4 Reproducibility of experiments

In order to allow interested readers to reproduce our experiments but also for further research and comparison with other classification techniques, we made all the datasets, algorithm implementation and testing source codes available on a dedicated web page¹. It contains the complete source code of the testbed, including the source code of all the benchmarked systems. The databases are also available, however using each dataset requires to report proper credits (different references for each dataset are available in annex of this document).

11.3 RESULTS AND ANALYSIS

We present in this section the results of the large scale study. We start by giving a detailed dataset-wise view on the results (Section 11.3.1). We further show in this analysis the statistical superiority of the HV-MOTS classifier over other schemes. We then try to obtain a higher-level view of the results by analyzing the potential influence of the number of features, samples, classes and domains on this statistical superiority (Section 11.3.2).

11.3.1 Dataset-wise results

We start by providing the complete classification results separated over each dataset for the HV-MOTS framework and compare it to the 1-NN, 5-NN, NC, SVM, MOTS and NP-MOTS classifiers. We provide the results based on overall classification accuracy and then draw our conclusions only from the complete statistical significance analysis.

Overall accuracy

We present in Table 14 the results of various methods with both the best and mean classification accuracies achieved for every methods. This figure exhibits the results for all the datasets studied.

[+ Critical difference graphs & analysis on Mean and on Best]

¹ <http://repmus.ircam.fr/esling/hvmots-datasets.html>

[+ Tukey-Kramer HSD on the best and mean results (solely)]
 [+ Same analyses on the `_COMPLETE_` set of features results (ALL features combos for ALL datasets)]
 [OVERALL ANALYSIS]
 [DATASET-WISE ANALYSIS]
 [BEST DIFFERENCE ANALYSIS]
 [WORST DIFFERENCE ANALYSIS]
 [MEAN / BEST ANALYSIS]

Statistical significance

As we discussed earlier, the reported classification accuracies are largely insufficient to draw any solid conclusion when comparing methods Demsar [105], Salzberg [337]. Therefore, we present in Figure 56 the results of statistical significance tests based on the mean column ranks obtained from Tukey-Kramer HSD based on a Friedman's ANOVA, as proposed by Demsar [105]. These statistical tests are performed on the accuracy matrix of every features combination available for each dataset.

[OVERALL ANALYSIS]
 [DATASET-WISE ANALYSIS]
 [BEST DIFFERENCE ANALYSIS]
 [WORST DIFFERENCE ANALYSIS]
 [MEAN / BEST ANALYSIS]

We further present in Figure 57 the statistical difference in mean classification accuracy from a one-way ANOVA method over every datasets. This allows to see the mean accuracy that can be expected by different methods for any combinations of features.

[OVERALL ANALYSIS]
 [DATASET-WISE ANALYSIS]
 [BEST DIFFERENCE ANALYSIS]
 [WORST DIFFERENCE ANALYSIS]
 [MEAN / BEST ANALYSIS]

11.3.2 *Global scale analysis*

The analysis of classification results in the previous section exhibited the strong statistical superiority of the HV-MOTS classifier. We now try to extract some higher-level clues on the statistical significance. We focus our attention on the different properties of the datasets, in order to see if the previous results can be explained by the variation in the datasets characteristics. Hence, we try to provide the eventual relationships between the characteristics of datasets (number of features, classes and samples) and the classification results. In order to perform this analysis, all the following presented results are, therefore, the mean column ranks of the Tukey-Kramer HSD over a Friedman ANOVA. Whenever we analyze a different aspect of variability, we concatenate all datasets that fall under a particular category and then compute the entire ANOVA. For example, in the *class-wise* analysis, the results of all 2-class datasets are concatenated in order to obtain a complete data matrix to analyze, and so on with higher number of classes. We perform this analysis for *class* cardinality, *samples* cardinality, *features* cardinality and finally regroup results depending on the scientific *domain* of study.

Name	1-NN	5-NN	NC	SVM	MOTS		
					-	NP	HV
Arabic digit	99.98	99.98	99.98	99.98	99.98	99.98	99.98
Artificial characters							
Australian signs							
Australian signs (HQ)							
BciIII-01-Tubingen							
BciIII-02-Albany							
BciIII-03a-Graz							
BciIII-03b-Graz							
BciIII-04a-Berlin							
BciIV-01-Berlin							
BciIV-03-Freiburg							
Biomag-2010							
Brain-Computer							
Challenge-2011							
Character-trajectories							
Dachstein							
Digit-hands							
Eeg-alcoholism							
Eeg-epfl							
Forte							
Gaitpdb							
Handwritten							
Ionosphere							
Japanese-vowels							
Libras							
Meg mind reading							
Neuro-emotional							
Neuro-visual							
Pen Characters							
Pen Characters							
Person activity							
Physical action							
Ptbdb							
Robot failures							
Slpdb							
Sonar							
Synemp							
Vfdb							
Vicon physical							
Wall-robot							

Table 14: Comparison of overall classification accuracies for different methods

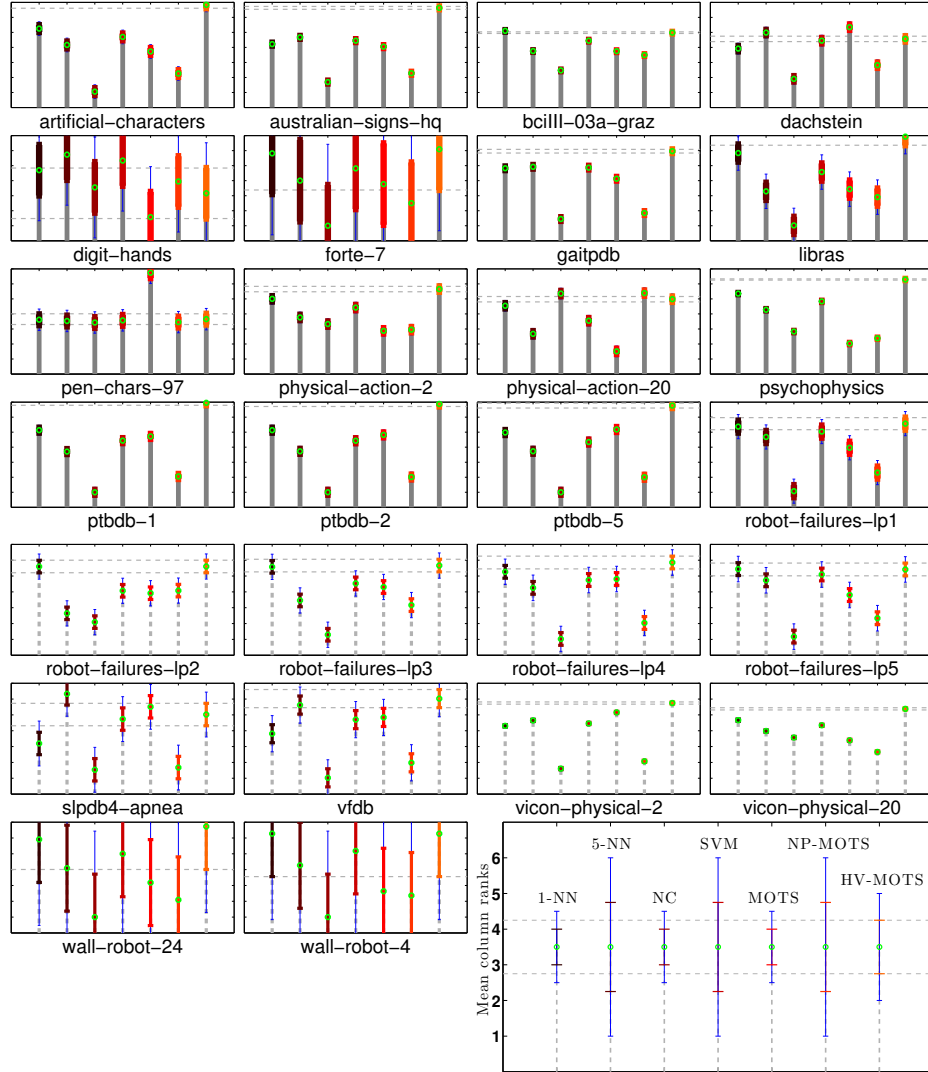


Figure 56: Comparison of statistical significance between classification methods based on the Tukey-Kramer HSD over Friedman's ANOVA.

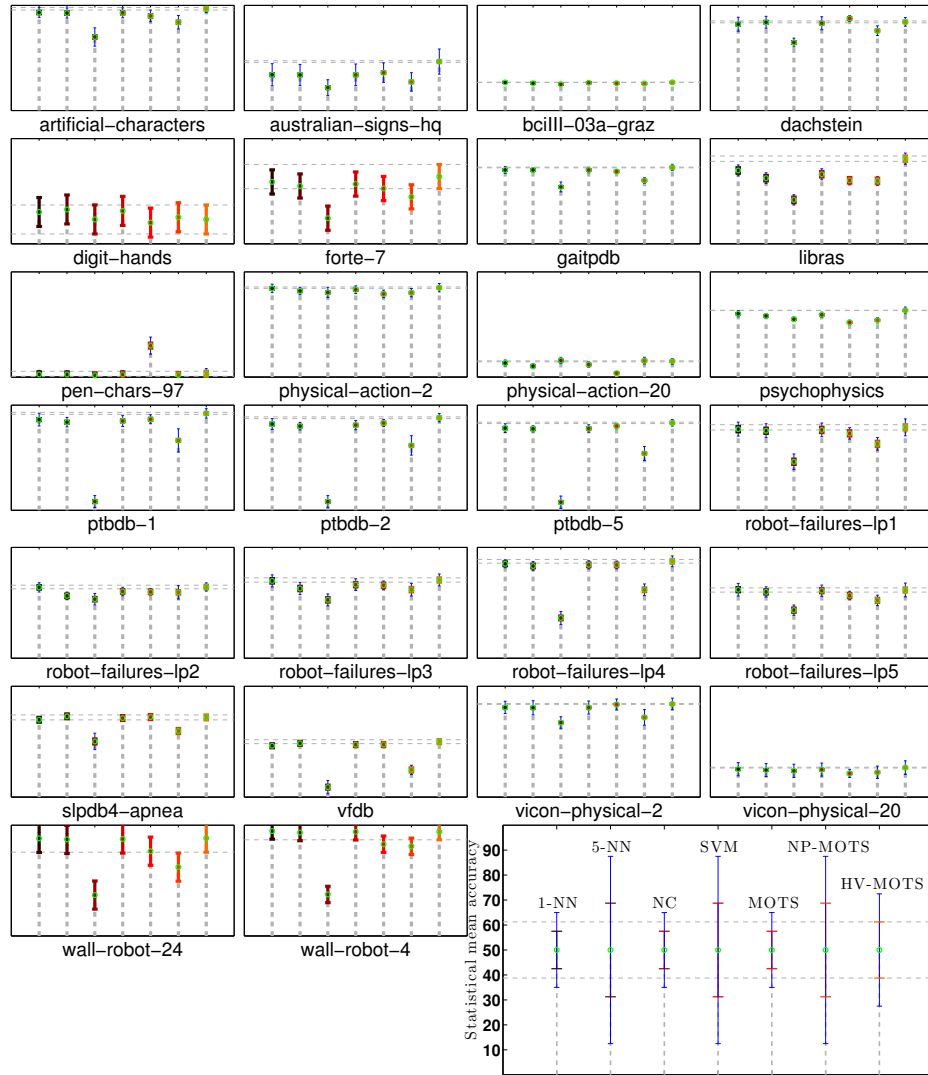


Figure 57: Comparison of statistical significance between classification methods based on the statistical mean difference over a one-way ANOVA.

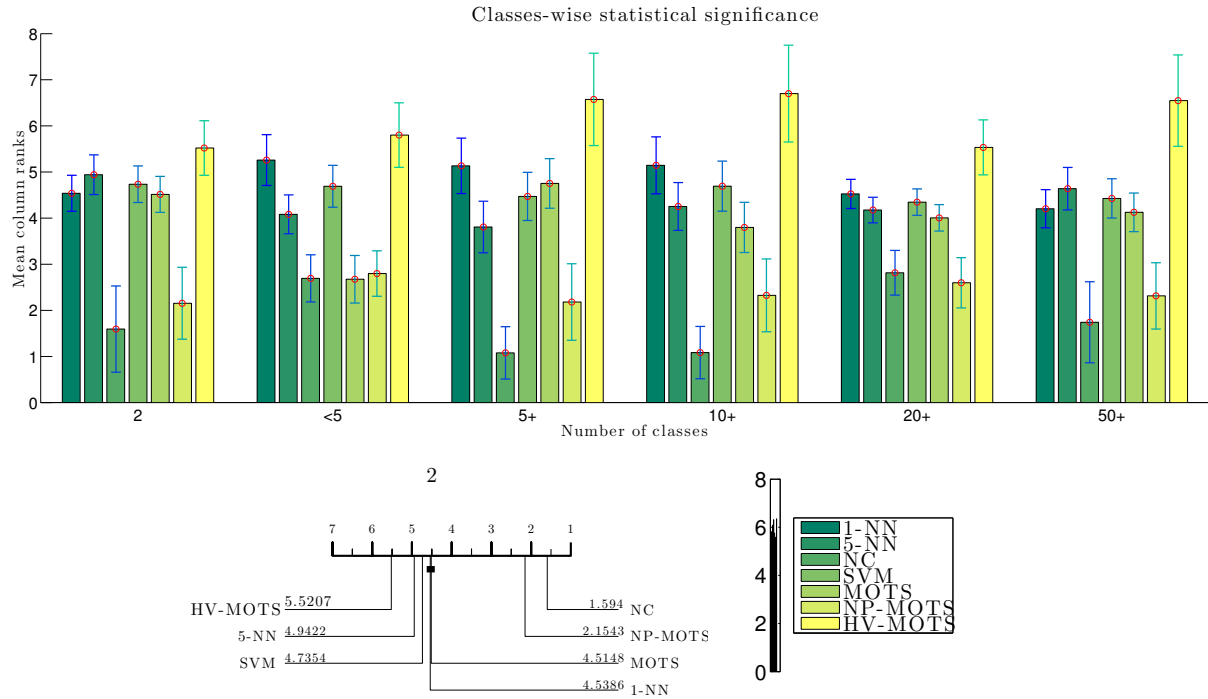


Figure 58: Comparison of statistical significance between classification methods for an increasing number of classes based on the Tukey-Kramer HSD over Friedman's ANOVA

Classes cardinality analysis

We first analyze the statistical difference between methods depending on the number of classes in the dataset. Therefore, we regroup classification results from datasets that share the same number of classes. We then perform the statistical significance analysis on various results matrix. The results are presented in Figure 58.

[OVERALL ANALYSIS / METHODS COMPARISON]

[INCREASING NUMBER ANALYSIS]

[BEST PERFORMANCE]

[WORST PERFORMANCE]

Samples cardinality analysis

We now analyze the statistical difference between methods depending on the number of samples in the dataset. Figure 59 shows the results of statistical differences between methods for an increasing number of samples.

[OVERALL ANALYSIS / METHODS COMPARISON]

[INCREASING NUMBER ANALYSIS]

[BEST PERFORMANCE]

[WORST PERFORMANCE]

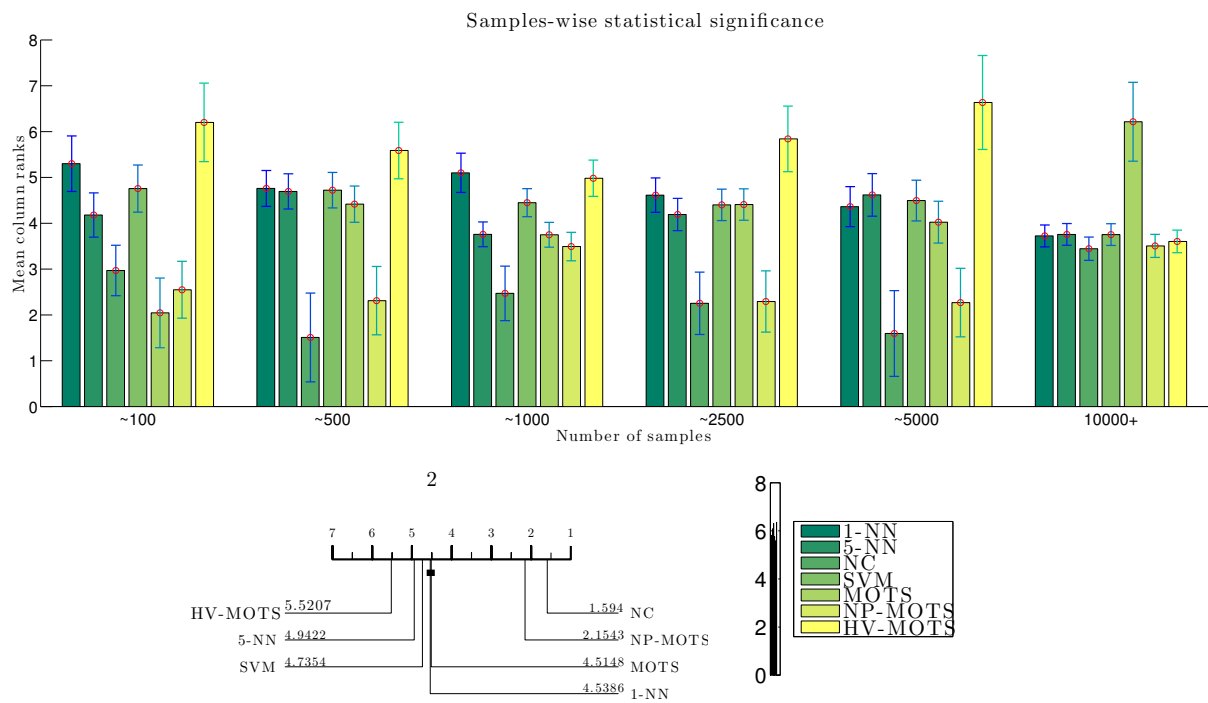


Figure 59: Comparison of statistical significance between classification methods for an increasing number of samples based on the Tukey-Kramer HSD over Friedman's ANOVA

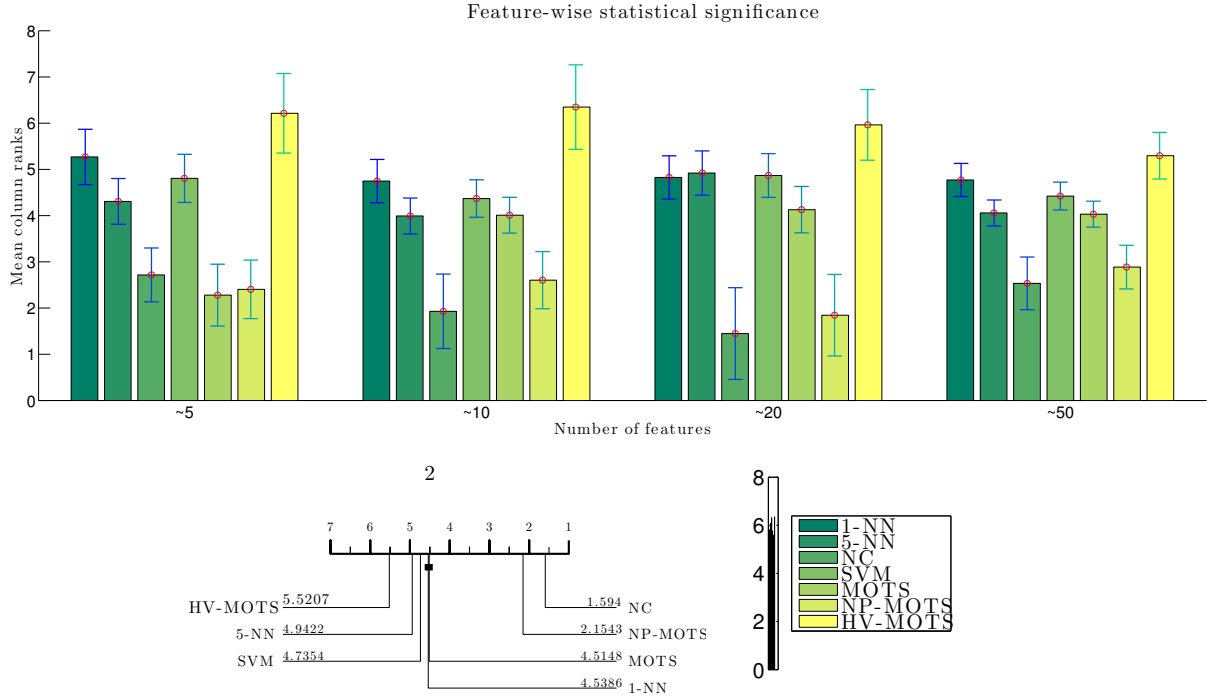


Figure 60: Comparison of statistical significance between classification methods for an increasing number of features based on the Tukey-Kramer HSD over Friedman's ANOVA

Features cardinality analysis

We now analyze the statistical difference between methods depending on the number of features in the dataset. Figure 59 shows the results of statistical differences between methods for an increasing number of features.

[OVERALL ANALYSIS / METHODS COMPARISON]

[INCREASING NUMBER ANALYSIS]

[BEST PERFORMANCE]

[WORST PERFORMANCE]

Domain-wise analysis

We now perform a domain-specific analysis in order to regroup dataset depending on their topics. Hence, we analyze if there is a link between the success of the HV-MOTS paradigm and the nature of underlying data. Figure 61 summarizes these results.

11.4 COMPARISON TO STATE-OF-ART RESULTS

We now provide a comparison of the HV-MOTS results to reported state-of-art accuracies and methods. This comparison is based on the classification accuracy reported in the corresponding research paper. We also provide the difference between the 1-NN

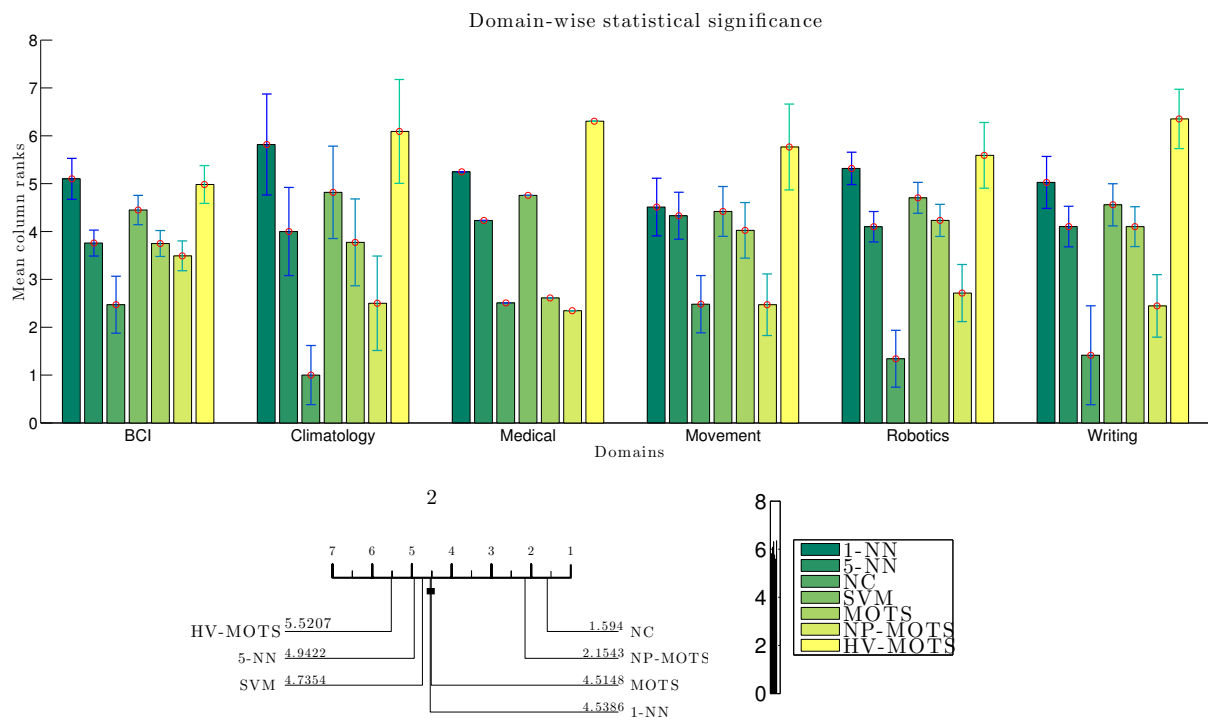


Figure 61: Comparison of statistical significance between classification methods depending on the scientific domain being studied based on the Tukey-Kramer HSD over Friedman's ANOVA

best classification accuracy. We recall that no feature pre-processing, extraction or modification have been performed whatsoever. Therefore, these results are based solely on the available raw time series and corresponding mean and deviation features. Table 15 summarizes the comparison to state-of-art algorithms on each of the dataset. We also provide in this table the reference papers along with the classification method used and the corresponding reported accuracy.

The first striking observation is the wealth and diversity of classification methods that have been used for these datasets. Furthermore, each of these researches have applied different and complex pre-processing steps to the features. Usually, the corresponding methods also undergo a very careful and extensive parameter tuning.

11.5 EXTENDED ANALYSIS

We provide in this section an extended analysis of more specific aspects of the results over the variety of datasets. First, we provide a detailed analysis of the features automatically selected to obtain the best classification results for each dataset. This allow us to get a deeper understanding on the reason *why* the classification is accurate. Then, we provide a deeper analysis on the influence of different parameters over the time series matching that could influence the classification results. Finally, we show the influence of using the time series data compared to mean and deviation statistics.

11.5.1 *Selected features*

The advantage of the proposed evaluation methodology is that it gives us access to the relative discriminative power of every feature available for each dataset. Therefore, we can easily investigate the features automatically selected by the testbed to obtain the best classification accuracy for each dataset. We can thus draw some interesting lines of future work for understanding the best features selected automatically.

11.5.2 *Warping or resampling*

We analyze in this section the influence of the warping window.

11.5.3 *The power of time*

Our approach is based on the assumption that the use of complete temporal information through the comparison of time series allows a better recognition and therefore a more accurate classification. Therefore, to validate this hypothesis, we present in figure ?? the dataset-wise analysis of results depending on the set of information used. Therefore, this figure contains the classification results without time series features (*mean* sets of information), with time series only (*temporal* sets without any mean or deviation information) and a combination of both sources (*mixed* sets).

Name	Algorithm	Results	HV-MOTS
Arabic digit	Vector Quantization + Tree	93.12%	100%
Artificial characters	Genetic Algorithm	98.68%	
Australian signs	Hidden Markov Model	71.2%	
Australian signs HQ	Tree Class Algorithm	94.5%	
BciIII-01-Tubingen	Common Spatial Subspace	91.0%	
BciIII-02-Albany	Gaussian SVM	73.5%	
BciIII-03a-Graz	Multi-Class CSP - Fisher ratios	0.7926	
BciIII-03b-Graz	Probabilistic Morlet Wavelets	89.3%	
BciIII-04a-Berlin	Neural Network - CSSD	92.98%	
BciIV-01-Berlin	Principal Component Analysis	0.382	
BciIV-03-Freiburg	SVM + LDA	46.9%	
Biomag-2010	SVM	69%	
Brain-Computer	Adaptive Auto-Regress+ SVM	85.75%	
Challenge-2011	Matrix of regularity	85.9%	
Character-trajectory	HMM + GMM	93.67%	
Dachstein	-	-	
Digit-hands	3-NN	97.8%	
Eeg-alcoholism	Multivariate HMM	78.5%	
Eeg-epfl	Bayesian LDA	95%	
Forte	Shared-NN	77.5%	
Gaitpdb	Neural Network - Wavelets	77.33%	
Handwritten	HMM + SVM	92%	
Ionosphere	Genetic Programming	94.2%	
Japanese-vowels	5-state continuous HMM	96.2%	
Libras	Spiking Neural Network	88.59%	
Meg mind reading	Bayesian Cross-Correlation	68.0%	
Neuro-emotional	-	-	
Neuro-visual	-	-	
Pen Characters	DTW + NN	89.15%	
Pen Characters	Template matching + NN	91.8%	
Person activity	Meta-Prediction Agents	91.33%	
Physical action	Genetic Programming	73.3%	
Ptbdb	Random Forest	75.1%	
Robot failures	Feature Transform + NN	80%	
Slpdb	Multi-Scale SVM	88.97%	
Sonar	Minimum Message Tree	76%	
Synemp	-	-	
Vfdb	Filter + Peak detection	91.5%	
Vicon physical	Dynamic Neural Network	95.4%	
Wall-robot	Polynomial SVM	95.58%	

Table 15: Comparison of classification accuracies with state-of-art results on the same datasets. We provide for each dataset the original algorithm used to obtained the reported classification accuracy.

Name	#	%	Features
Arabic digit			
Artificial characters			
Australian signs			
Australian signs (HQ)			
BciIII-01-Tubingen			
BciIII-02-Albany			
BciIII-03a-Graz			
BciIII-03b-Graz			
BciIII-04a-Berlin			
BciIV-01-Berlin			
BciIV-03-Freiburg			
Biomag-2010			
Brain-Computer			
Challenge-2011			
Character-trajectories			
Dachstein			
Digit-hands			
Eeg-alcoholism			
Eeg-epfl			
Forte			
Gaitpdb			
Handwritten			
Ionosphere			
Japanese-vowels			
Libras			
Meg mind reading			
Neuro-emotional			
Neuro-visual			
Pen Characters			
Pen Characters			
Person activity			
Physical action			
Ptbdb			
Robot failures			
Slpdb			
Sonar			
Synemp			
Vfdb			
Vicon physical			
Wall-robot			

Table 16: Comparison of classification accuracies with state-of-art results on the same datasets.
144 We provide for each dataset the original algorithm used to obtained the reported classification accuracy.

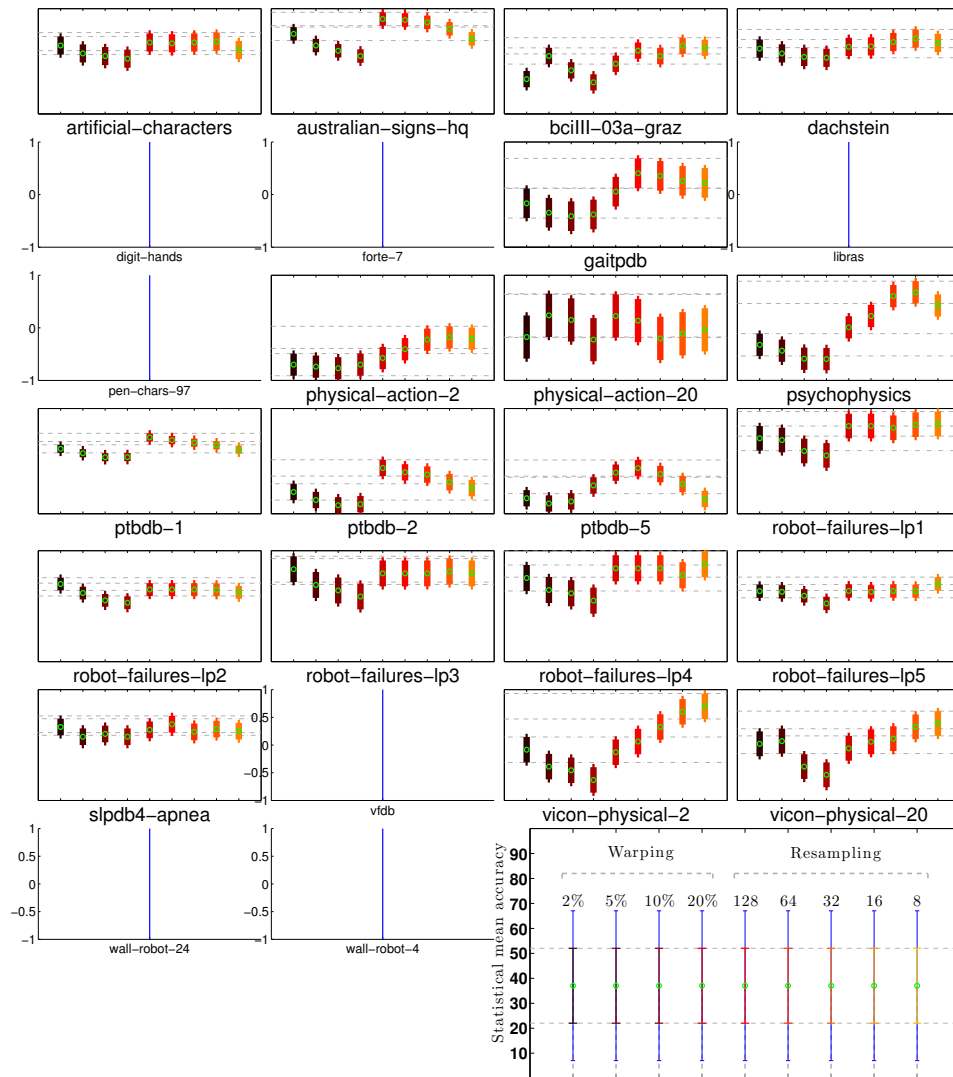


Figure 62: Comparison of the statistical significance for different parameters of warping for the DTW distance as opposed to simple resampling factors of the time series compared with the Euclidean distance.

12

UNICITY OF HEART SOUNDS

As every organ, the human heart follows a physiological development singular to each living being. Its formation along the gestation steps of prenatal development makes it the first functional organ in an embryo. We show that, based on the HV-MOTS classifier, we can construct the first system that accurately identifies someone through the sounds its heart produces. We attain this goal by considering listening as an art and finding inspiration in musical analysis. This system is able to attain error rates equivalent to established biometric traits such as speech or gait recognition. We show *how* to listen to the peculiarities of each heart by developing a specific set of features based on the Stockwell transform, called the *S-Features*. Our findings are supported by the largest PCG dataset ever collected. This includes the Mars500 isolation study of the Russian, Chinese and European spatial agencies. This dataset contains heart sound recordings that were recurrently collected over a time span of almost two years. The complete set of data allows to support biometric identification over large numbers of persons, long time spans and different physiological states. We also provide the first study ever on the phenomenon of template ageing for heart signals.

12.1 BIOMETRIC SYSTEMS

Throughout its complete cycle, the heart produces a characteristic sound signature. Two major components usually emerge from these cardiac contractions. Listening to these sounds for the purpose of detecting anomalies in the circulatory system has been a common practice for the past two centuries. Hence, medical auscultation is taught as a form of art, which implies listening to the music of hearts with a firm clinical knowledge. Given this tradition, we could wonder if the sounds produced by a heart are utmostly *unique* to each individual. If we can exhibit this uniqueness through careful listening, we could be able to identify which person each heart belongs to with a single heart beat recording. Therefore, we study the characteristics of heart sounds as biometric features. Biometric systems offer a natural and secure solution to authentication paradigms [190]. These methods seek to identify a person based on its distinguishing physiological or behavioral characteristics. The core strength of such systems is that the feature used for authentication is derived directly from the user. Therefore, it is unlikely to be lost, forged or stolen as it might be the case for *token-based authentication* (keys, passwords or cards). A biometric system is essentially a *pattern matching* method. Hence, it establishes the authenticity of an identifying characteristic, by comparing it to a template collected in an enrollment procedure. Humans have always subconsciously applied such biometric recognition principles by analyzing the characteristics of the face, voice or even gait in order to identify other human beings. The most commonly known biometrics, like fingerprints or DNA, are now being used as wide spread international standards for identification with the emergence of biometric passports. Nevertheless, emerging biometrics are also being increasingly studied like facial thermogram, retinal scan, vocal prints or hand geometry. An ideal biometric feature should be *universal* throughout the population, *unique* to each person, *permanent* over time and easily *collectable* [190]. The biometric system

itself should exhibit good *performance*, have a good *acceptability* by the population and prevent fraudulent attacks by *circumvention* [191]. Only very recently has emerged the idea of using the electrical activity of the heart (*ElectroCardioGram* (ECG)) as a biometric feature [47]. This information has been studied either through the direct extraction of representative points (called *fiducial points*) [366] or through spectral analysis [397]. Despite the nontriviality of using this information for identification, it has recently been shown that the features extracted from an ECG are invariant to the sensor location and physiological state of the subject [188]. The idea of using heart sounds should follow as a logical consequence [301]. However, studies in that line of research have been extremely limited and exhibited several flaws. First, most of the studies are evaluated on very narrow datasets, usually less than twenty subjects [121, 195, 301]. Some studies even make use of the same recording for enrollment and identification [302]. These methodologies cannot truly support any finding in a biometric context. Finally, the few studies that use larger databases with distinct recordings exhibit poor performances [38]. We provide in annex (cf. Section A.2.1) a complete comparison of existing biometric studies.

We show that we can access to the peculiarities of heart sounds if we consider listening as an art form and know *where*, *when*, *what* and *how* to listen. To know *where* we should listen, we study the unique frequency distribution of heart beats over a wide scope of subjects. To know *when* to listen, we use a segmentation procedure in order to focus our listening on each heartbeat separately. To know *what* to listen, we develop a specifically-tailored set of features based on the Stockwell Transform and inspired by *Music Information Retrieval* (MIR). Finally, to know *how* to listen we use the HV-MOTS classification paradigm.

12.2 HEART PHYSIOLOGY

The human heart is the muscular pump that manage the blood circulation throughout every parts of our body. The mechanisms and actions of the heart are distributed over its left and right halves. The right part dispatches the blood to the lungs and therefore dispenses pulmonary circulation. The left heart provides the supply of oxygen and nutrients to our entire body. Each half is further divided into two cavities known as the *atrium* and *ventricle*, where a set of valves (*tricuspid*, *mitral*, *pulmonary* and *aortic*) regulates our blood flow. Each of these components contribute in a specific way to the complete cardiac cycle along the periods of muscle contraction (*systole*) and relaxation (*diastole*). Throughout its complete cycle (illustrated in Figure 63), the heart produces a characteristic sound signature constituted of two major components. There is no consensus on the precise physiological origins of every of their acoustic components. Some potential origins might be valvular (atrioventricular and exhaust valves), muscular (ventricular muscle) and circulatory (intracardiac blood flow) actions [256, 270, 377]. Even within this mixity, some major contributors have been clearly established through simultaneous PCG and echocardiography recordings. The first sound (S_1) partly originates from the closure of the mitral valve and the subsequent (8 to 12 milliseconds later) opening of the aortic valve at the end of the isovolumic contraction. The second sound (S_2) is induced by the closure of the aortic valve after the period of isovolumic relaxation. The onset of the S_2 sound coincide the end of systolic contraction and the occurrence of the pulmonary ejection. The remaining period of diastolic rest is normally not audible in healthy subjects.

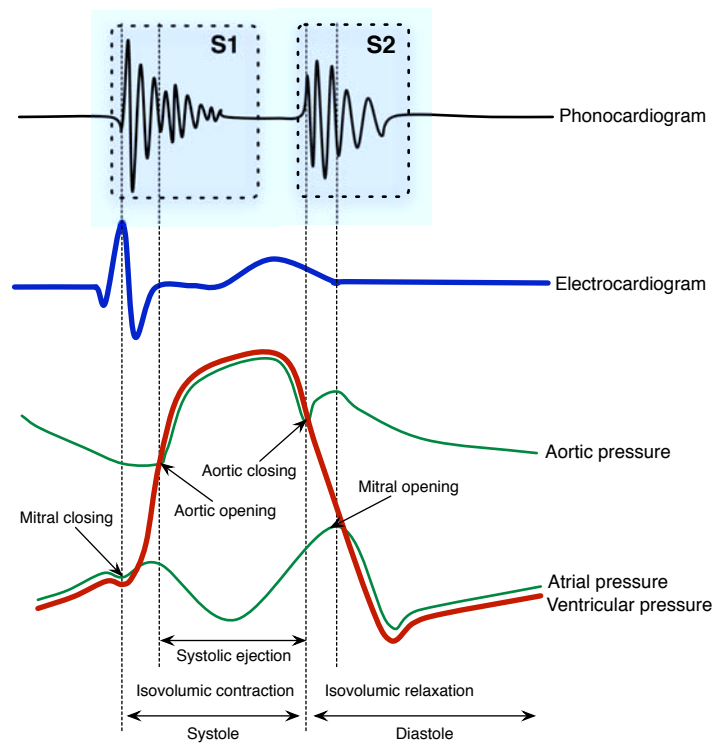


Figure 63: A complete cardiac cycle analyzed through recordings of the heart sounds (PCG), its skin electrical activity (ECG) and pressure in the aortic and atrial valves.

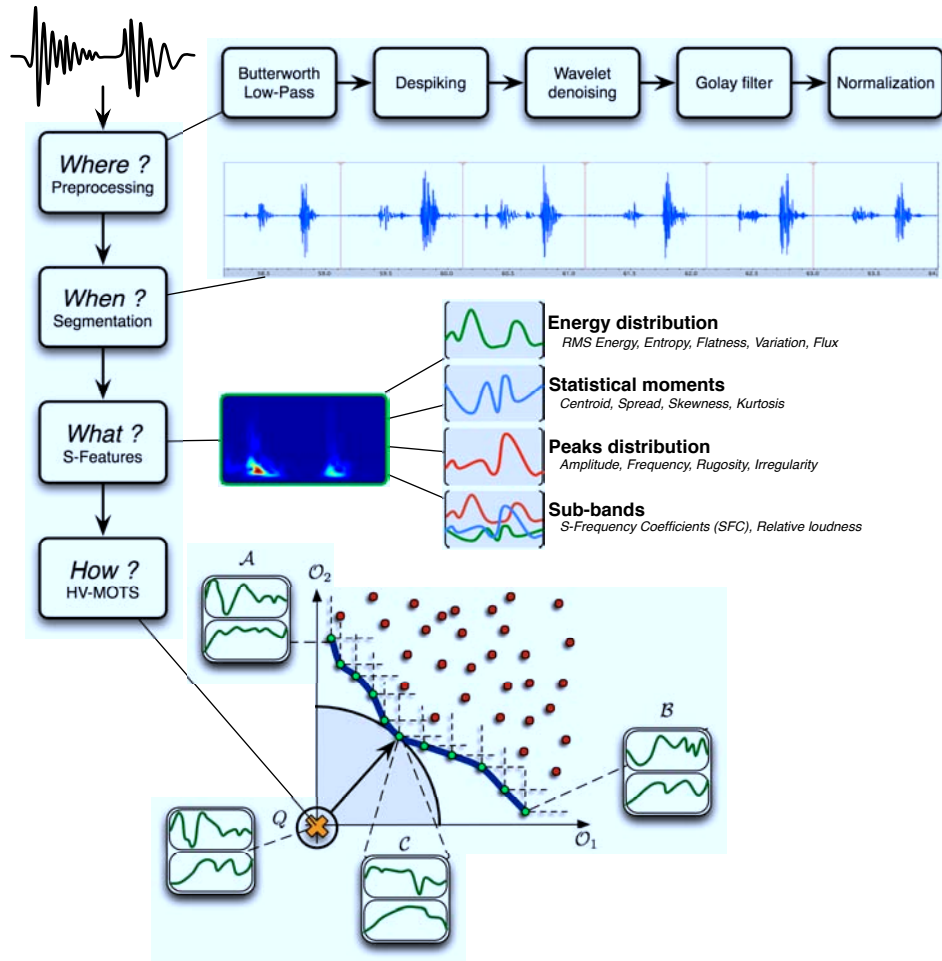


Figure 64: Algorithmic workflow for our heart sounds biometry system, summarizing the four milestone of listening which are *where*, *when*, *what* and *how* to listen.

12.3 LISTENING TO THE HEART

Cardiac sounds are cyclic and therefore repetitive phenomena. A single cardiac cycle can thus be considered as the elementary unit of this study. We will therefore focus on listening from the beginning of systole, to the end of diastole. We will show that this period is consistent between several cycles of an individual by using the workflow displayed in Figure 64.

12.3.1 Where to listen (pre-processing)

As heart sounds possess highly specific characteristics, they require very cautious analysis. First, we need to uncover the useful bandwidth of cardiac sounds. Hence, we will be able to focus our listening on the range where the interesting spectral information lies. Figure 65 displays the typical frequency distribution of human heart sounds. This distribution is the normalized energy of the spectral information for 15,814 complete cardiac cycles recorded from 212 different subjects. We can see that the cardiac sound information primarily lies between 5 to 400 Hz, which makes it an

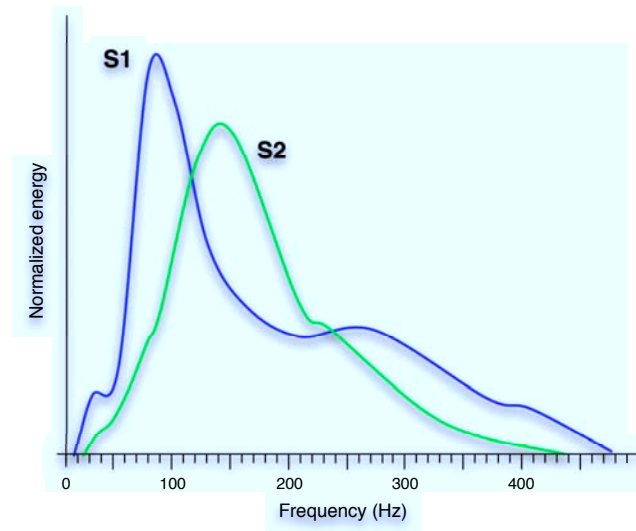


Figure 65: Typical frequency distribution of heart sounds based on the Stockwell transform analysis of 15,814 complete cardiac cycles. The S1 and S2 sounds recorded from 212 different subjects have been processed separately.

significantly low-frequency signal. Therefore, we first process the heart signals with a [5, 500]Hz band-pass filter, in order to remove all non-relevant information. As the present study is centered on the *sounds* produced by human hearts, it is subject to the same constraints as any audio-related system. Therefore, we need to remove the artifacts caused by noisy and low-quality recording conditions. These shortcomings are usually encountered due to the levels of surrounding noise but also the jitters caused by defects in the recording device itself. In order to alleviate these problems, we first process the signal with a despiking algorithm based on a phase space decomposition. Then, we apply a wavelet denoising algorithm with the 6th Daubechies wavelet (because of its similarity with the PCG signal). We can expect widely varying loudness levels in the resulting signals. Therefore, we apply a Golay filter of degree 9 with a window of 65 samples and an amplitude normalization procedure. This allows to enhance the signal of each heart beat. We provide in annex of this document (cf. Section A.3) a complete analysis of the impact from each pre-processing component on the overall performances.

12.3.2 When to listen (*segmentation*)

As we intend to use single cardiac cycles as elementary units, we need to precisely determine the boundaries of each heart beat. Therefore, we will focus on listening from the beginning of systole, to the end of diastole. Once these segments have been detected, we can extract independent cycles for subsequent analysis. Therefore, we perform a transient segmentation based on the difference in spectral flux between successive frames. An energy variation threshold then allows to filter less significant transients. We apply an analysis window of 0.7 seconds with an 8 times oversampling. We extract and normalize each cardiac cycle and then store the resulting signals in separate files. Therefore, the identification template for each individual is provided by the complete set of cardiac cycles extracted from a single recording. We also provide in annex of this

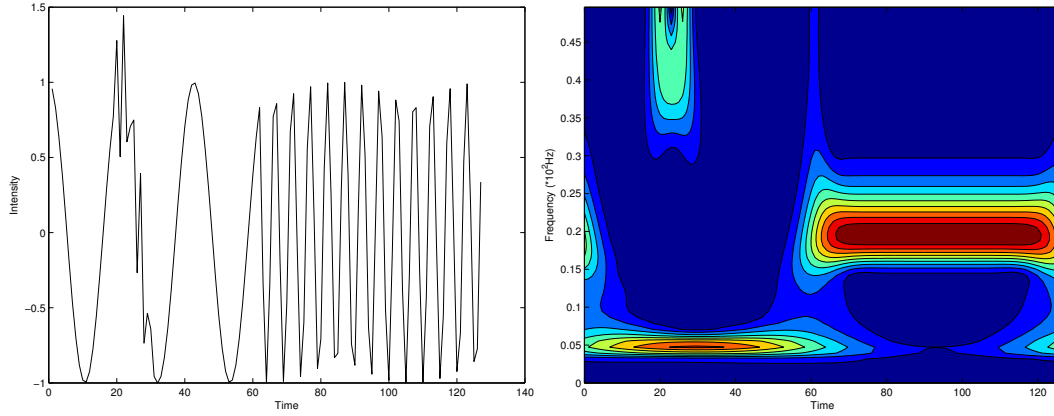


Figure 66: Illustration of the temporal and frequency resolution of the Stockwell transform. A synthetic signal (left) with very specific properties and its corresponding S-Transform spectrum (right).

document (Section A.3) an extensive analysis of the influence of each segmentation parameter on the performance of the system.

12.3.3 What to listen (*S-Features*)

The Stockwell transform

The considerably low-frequency properties of heart sounds require a spectral transform that could adapt to these unique characteristics. Various signal processing tools have been developed to access the spectral information, such as the Short-Time Fourier Transform (STFT) or Continuous Wavelet Transform (CWT). Unfortunately, heart sounds have an extremely narrow bandwidth with most of their energy radiating below the frequency resolution of these traditional decompositions. The Stockwell transform (or *S-Transform*), originally developed for analyzing geophysical data [361] provides an adequate solution to this problem. It is defined as a generalization of both the STFT and the CWT and overcomes some of their limitations. The S-transform of a function $h(t)$ can be defined as a CWT with a gaussian mother wavelet multiplied by a phase factor [361]

$$S_x(t, f) = \int_{-\infty}^{\infty} h(\tau) \frac{|f|}{\sqrt{2\pi}} e^{-\frac{(t-\tau)^2 f^2}{2}} e^{-i2\pi f \tau} d\tau \quad (12.1)$$

The S-Transform exhibits a frequency-dependent resolution. This leads to an extremely fine resolution even at very low frequencies (below 50Hz) where lies the cardiac bandwidth. Therefore, it allows a better distinction of spectral components relevant to the cardiac information. Furthermore, unlike the CWT, modulation sinusoids are fixed with respect to the time axis. This property localizes dilations and translations and thus provides the same temporal resolution for every frequency bins. Figure 66 illustrates these properties. We use the fast S-Transform algorithm proposed in [61] which strongly reduces its computational complexity.

S-Features

The S-Transform provides an highly precise representation of the temporal evolution of frequency distributions for each cardiac cycle. However, the properties singular to each heart remain strongly diluted in such an overwhelming quantity of information. In order to uncover a potential uniqueness, we need to focus our listening by extracting higher-level information. To that end, we developed a specifically-tailored set of high-level features (called *S-Features*) based on recent research in music analysis (MIR). For the sake of clarity, we only provide here a summary of the implemented features but the complete mathematical definitions are available in annex of this document (cf. Section 12.3.3). First, we study the evolution of the S-Transform distributions by computing their statistical moments. Therefore, we calculate the mean (*centroid*), variance (*spread*), symmetry (*skewness*) and peakedness (*kurtosis*) of the successive frequency distributions. We further study the evolution of these distributions through their *energy*, *entropy*, *brightness*, *flatness*, *rolloff* and *variation*. We adapt the Mel-Frequency Cepstral Coefficients (MFCCs) for describing the evolution of spectral shape over various frequency bands through *S-Frequency Coefficients* (SFCs) (we detail this feature in the next section). Finally we study the evolution of peaks in the distribution by computing the relative differences in their frequencies (*roughness*) and amplitudes (*irregularity*). We resample each feature to a fixed length which allows us to abstract from the heart rate and therefore accounts for various physiological conditions.

S-Frequency Coefficients (SFC)

The S-Frequency Coefficients (SFCs) provide a description of the evolution of spectral shape over various frequency bands. Therefore, the computation of the SFCs follows a computing scheme analogous to the MFCCs. However, the MFCCs process the spectral components through a filterbank based on the Mel scale in order to approximate perceptual results in pitch judgements. However, in the case of heart sounds, the Mel scale is not as relevant. Therefore, we use a logarithmic filterbank with specific distributions for the filters center frequency. If we look at the distribution of heart sounds (figure 65), we can see that most of the energy is concentrated under 200Hz. The overall bandwidth lies between 10Hz and 500Hz. Therefore, we compute the frequency centers

$$f_{\text{center}}^i = (f^{\text{N}_{\text{bands}}}/f^1)^{\frac{i-1}{\text{N}_{\text{bands}}-1}} \cdot f^1$$

given a variable number of bands N_{bands} and $f^1 = 10$ the center of the first filter and f^{N} the center of the last filter. Note that this design will lead to extremely narrow filters in the low frequency range. Using such filters is possible with the S-Transform which provides one bin per frequency value. Finally, a Discrete Cosine Transform (DCT) is used given its energy compaction property to reduce the dimensionality. We also compute the first derivative (DSFC) and second derivative (DDSF) of the SFCs.

12.3.4 How to listen (HV-MOTS Scoring)

As we will discuss later (Section 12.4.2), the input evaluation required by biometrics system is quite different from classification tasks. Indeed, an identification attempt might be rejected if no template in the database match the input. Therefore, a biometric system should be able to make such a decision. To provide this subtlety, we maintain

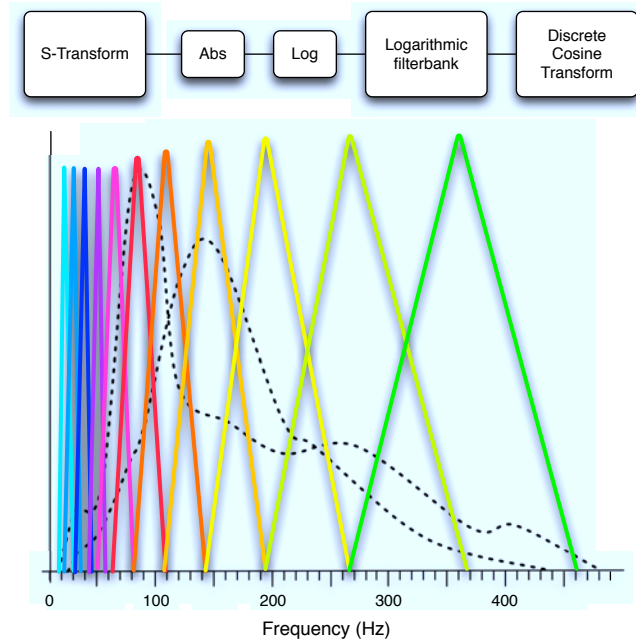


Figure 67: Computation workflow (up) and specific filter design (down) of the S-Frequency Coefficients (SFC). As we can see, the SFC are computed on a model similar to the MFCC. However, its fundamental differences comes from the use of the S-Transform and its filterbank designed to match the properties of heart sounds.

the workflow of HV-MOTS classification until the final decision. Hence, we do not select a particular class but rather keep the *normalized hypervolume* of each class as a matching score.

12.4 EXPERIMENTS

12.4.1 Datasets

We collected two datasets of PCG recordings, each of them having inherent advantages and flaws. Both datasets meet the main requirement to provide at least two distinct recordings for each person, collected at different times in separate sessions. First, the HSCT-11 dataset collected by Beritelli and Spadaccini [38] is the largest PCG set available of two distinct recordings per individual. This collection allows to assess the feasibility of heart sounds biometry over a large number of subjects. Nevertheless, the biggest flaw of this dataset is that it covers restricted time spans, most of the recordings being collected on the same day a few hours apart. Moreover, these recordings were collected from subjects in resting states only, which obviously can not account for varying physiological conditions. To overcome these shortcomings, we also include in our study a unique PCG dataset collected throughout the Mars 500 isolation experiment. This study is a joint Russian, European and Chinese spatial agencies psycho-social experiment. It was designed to analyze the psychological and physiological effects of a long-term deep space mission. Hence, it reproduces the conditions of a complete manned roundtrip to Mars. This dataset is the first of its kind for cardiac studies, and its advantages in our context are manifold. First, its uniqueness stems from

	Gender		Ages			Records separation			Physical state
	M	F	Min	Max	Med	Min	Max	Mean	
HSCT-11	157	49	15	96	28	1	2	1.1	Rest
Mars500	6	0	26	38	31	19	472	236.5	Mixed
Both	163	49	15	96	28	1	472	44.75	Mixed

Table 17: Details of the PCG recordings datasets. We provide the

the time spans over which auscultations have been performed. Indeed, PCG signals have been consistently recorded every three months over a complete period of almost two years of auscultation. This makes it the most detailed and longest systematical examination of individual heart sound recordings. Therefore, this dataset can unveil long-term evolutions and variabilities of heart sounds. Longer time intervals usually raise the difficulty in matching samples due to the phenomenon known as *template ageing*. This provokes an increase in error rates caused by time-related changes in the biometric pattern. We will therefore provide a specific evaluation framework designed to assess this effect. In addition, recordings have been performed under varying physical conditions which allows to account for fluctuating heart rates and diverse physiological factors. Moreover, recordings were collected by *auto-auscultation* from different non-experts in the medical field. This implies inherent *presentation* and *channel effects* with various levels of noise, placement of stethoscope and so forth. This also allows to account for the usability of such an identification system by novice users. All these characteristics provide an interesting range of robustness issues on the nature and variability of heart sounds. However, the Mars 500 dataset is flawed by its cardinality of six volunteers only. Therefore, we combine both datasets to evaluate the performance of our method. We follow the guidelines of [259] by providing a detailed description of the datasets. Table 17 summarizes with the demographics (gender and age) of the volunteers, time separation between samples and physical states under which recordings are made. Both dataset were collected with *cooperative, non-habituated* and *private* users with an *overt* and *attended* capture by an *open* system in a *standard* environment.

HSCT-11

This database contains heart sounds acquired from 206 subjects (157 males and 49 females). The files were recorded in Wave format using a sampling frequency of 11025 Hz with 16 bits per sample. Two separate recordings were collected from each person (the length of the recordings varying from 20 to 70 seconds). Both recordings were usually collected the same day, separated by a break. Every person was sitting, in resting state while the auscultation was performed near the pulmonary valve.

Mars500

This dataset is part of the *CardioPsy* study conducted in the Mars500 study. The hardware used was a digital stethoscope developed by the INFRAL society. The sounds are recorded by a piezoelectric microphone with linear response between 20 and 4000 Hz, AD converted and transmitted with Bluetooth to a laptop as hosting device. There is no pre-filtering of the sounds which are recorded at an 8000Hz sampling frequency with 16 bits per sample in Wave format. The dataset contains 54 sounds from the

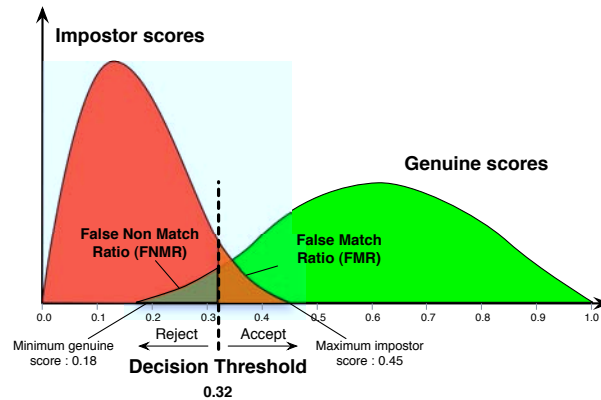


Figure 68: The evaluation of a biometric system in a real-life scenario given its distributions of *genuine* and *impostor* scores.

six volunteers, corresponding to nine cardiac auscultations that have been collected periodically over 520 days.

12.4.2 Evaluation methodology

Unlike classification problems where an input is bound to be labeled as belonging to one of the classes, a biometric system should also be able to detect if the incoming attempt does not belong to its template database. Therefore, the matching algorithm should decide whether the input is a genuine or impostor attempt given a set of individuals. To make this decision, the matching algorithms usually rely on a similarity score and a decision threshold (DT). If the score is inferior to the DT, the identification is negative (*impostor* attempt). Reciprocally, if the score is superior to the DT, the identification is positive (*genuine* attempt). Hence, biometric systems are prone to two kind of errors. First, when the algorithm confuses an impostor and identifies it as a genuine attempt from someone inside the database. This type of error is called a *False Match* (FM). Second, when the algorithm fails to identify a genuine attempt (ie. fails to find a person which is truly enrolled in the database). This type of error is called a *False Non-Match* (FNM). There is a tradeoff between these two kinds of errors depending on the setting of the DT. Raising the DT reduces the number of FM errors but increases the number of FNM errors and vice versa. These concepts are illustrated in Figure 68.

For the evaluation procedure, the first recording is used for computing the reference templates stored in the database. Every further recording is used to compute the set of matching scores. For each of the N persons in the database, the system computes one genuine score, and $N - 1$ impostor scores. This yields a final number of N genuine and $N \cdot (N - 1)$ impostor matching scores. As proposed by [151], we follow a thorough evaluation scheme and perform a multi-order analysis of results. This allows to obtain the most complete view on our results as possible.

12.4.3 Results

We present here the results of identification over a combination of both datasets. This first part allows to assess the feasibility of biometric identification over a wide range of

	Genuine		Impostor		D-Prime
	μ	σ	μ	σ	
HV-MOTS	0.987	0.014	0.721	0.131	1.979
MOTS	0.933	0.081	0.872	0.182	1.401
1-NN	0.895	0.016	0.782	0.142	1.326
5-NN	0.781	0.125	0.141	0.131	1.133

Table 18: Result of *Order-o analysis* for different methods

	1	2	3	4	5	6	7	8	9
Genuine									
μ	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
σ	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
Impostor									
μ	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
σ	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
D-Prime									

Table 19: Result of *Order-o analysis* for different levels using the HV-MOTS method

persons. We then try to understand the origins of such differences. Therefore, we further study the effect of very long time spans (or *template ageing*) thanks to the Mars500 dataset. Finally, we provide a comparison of our proposal to current state-of-art systems performance for various biometric features. We start by providing a multi-order analysis of results as proposed by several current researches on the assessment of biometric systems [151].

Order-o analysis

We first collect the *order-o statistics* by computing all possible genuine and impostor scores and count total number of FM and FNM errors to compute the cumulative measurements at fixed rates of the DT. The *False Match Rate* (FMR) gives the percentage of FM errors and the *False Non-Match Rate* (FNMR) gives the percentage of FNM errors for a given DT. Table 18 summarizes the results of the *Order-o analysis* between different methods.

Table 19 summarizes the results of the *Order-o analysis* for an increasing number of objectives using the HV-MOTS classifier.

Figure 69 shows the tradeoffs between the FMR and FNMR for different methods.

Order-1 analysis

We then compute the *order-1 statistics* by computing and plotting the *trade-off curves* which allows to show the performance of the system over a range of decision criteria. The *Detection Error Trade-off* (DET) curve represents the co-evolution of FMR vs. FNMR, by varying the DT. The curves are plotted parametrically on a log-log scale. The *Receiver Operator Characteristic* (ROC) curve, which is similar to DET curve, but

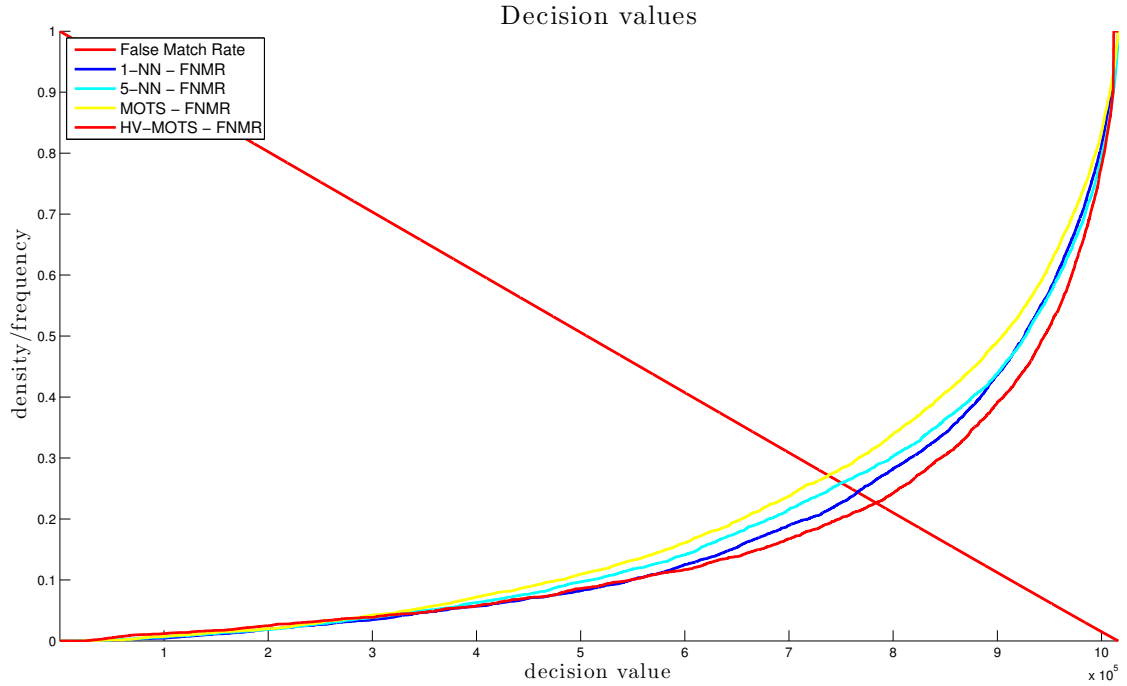


Figure 69: Possible tradeoffs between the False Match Rate (FMR) and the False Non Match Rate (FNMR) for different methods

plots the *True Acceptance Rate* ($TAR = 1 - FNMR$) against FMR. Table 20 summarizes the results of the *Order-1 analysis* between different methods.

Table 21 summarizes the results of the *Order-1 analysis* for an increasing number of objectives using the HV-MOTS classifier.

Figure 70 shows the *Detection Error Trade-off* (DET) curves for different methods.

Order-2 analysis

The *Order-2 analysis* computes the same statistics for all possible threshold values and computes the *Rank-1 identification rate* which is the number of times the correct person has the highest score of the database. We also compute the *Equal Error Rate* (EER) which is the percentage of errors induced by the DT setting that produces the same ratio of

False Non-Match Rate				
FMR =	0.1	0.01	0.001	0.0001
HV-MOTS	0.987	0.014	0.721	0.131
MOTS	0.933	0.081	0.872	0.182
1-NN	0.895	0.016	0.782	0.142
5-NN	0.781	0.125	0.141	0.131

Table 20: Result of *Order-1 analysis* for different methods

FNMR	1	2	3	4	5	6	7	8	9
0.1	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
0.01	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
0.001	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
0.0001	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013

Table 21: Result of *Order-1 analysis* for different levels using the HV-MOTS method

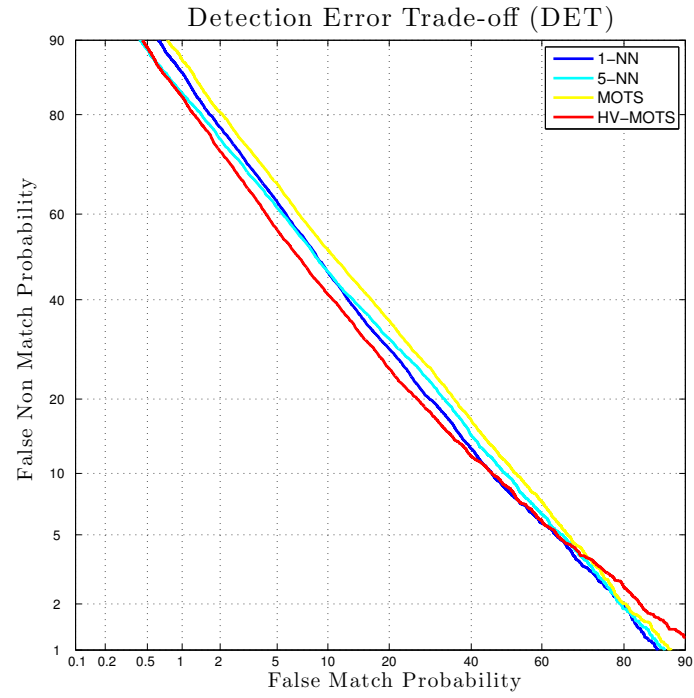


Figure 70: Detection Error Trade-off (DET) curves for different methods

	EER	Area Curve	Area Hull	Acc. (o)	Acc. (Max)
HV-MOTS	7.84	0.8891	0.9133	0.1345	0.9951
MOTS	11.23	0.7613	0.7311	0.0414	0.4116
1-NN	13.01	0.7821	0.9133	0.1345	0.9951
5-NN	14.13	0.6613	0.7311	0.0414	0.4116

Table 22: Result of *Order-2 analysis* for different methods

	1	2	3	4	5	6	7	8	9
EER	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
Area Curve	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
Area Hull	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
Acc. (o)	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
Acc. (Max)	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013

Table 23: Result of *Order-2 analysis* for different levels using the HV-MOTS method

FM and FNM errors. The EER provides a good summary of an algorithm performance on both types of errors. Table 22 summarizes the results of the *Order-2 analysis* between different methods.

Table 23 summarizes the results of the *Order-2 analysis* for an increasing number of objectives using the HV-MOTS classifier.

Figure 71 shows the *Receiver-Operator Characteristic* (ROC) curves.

Order-3 analysis

Finally we perform the *order-3 analysis* by plotting the *Cumulative Match Characteristic* (CMC) curve which displays the *Rank-k identification rates* for all possible ranks $k \in [1 \dots n]$. We then compute the distribution of *confidence intervals* for several performance indices. We further study the *scalability* of our approach by plotting these performance metrics over randomly selected increasing datasets and studying the corresponding trade-off curves. Table 24 summarizes the results of the *Order-3 analysis* between different methods.

Table 25 summarizes the results of the *Order-3 analysis* for an increasing number of objectives using the HV-MOTS classifier.

Figure 72 shows the *Rank-k identification rates* for all possible ranks.

	Rank-1	Rank-2	Rank-5	Rank-10	1 st 2 nd (μ)	1 st 2 nd (σ)
HV-MOTS	7.84	0.8891	0.9133	0.1345	0.9951	0.9951
MOTS	11.23	0.7613	0.7311	0.0414	0.4116	0.4116
1-NN	13.01	0.7821	0.9133	0.1345	0.9951	0.9951
5-NN	14.13	0.6613	0.7311	0.0414	0.4116	0.4116

Table 24: Result of *Order-3 analysis* for different methods

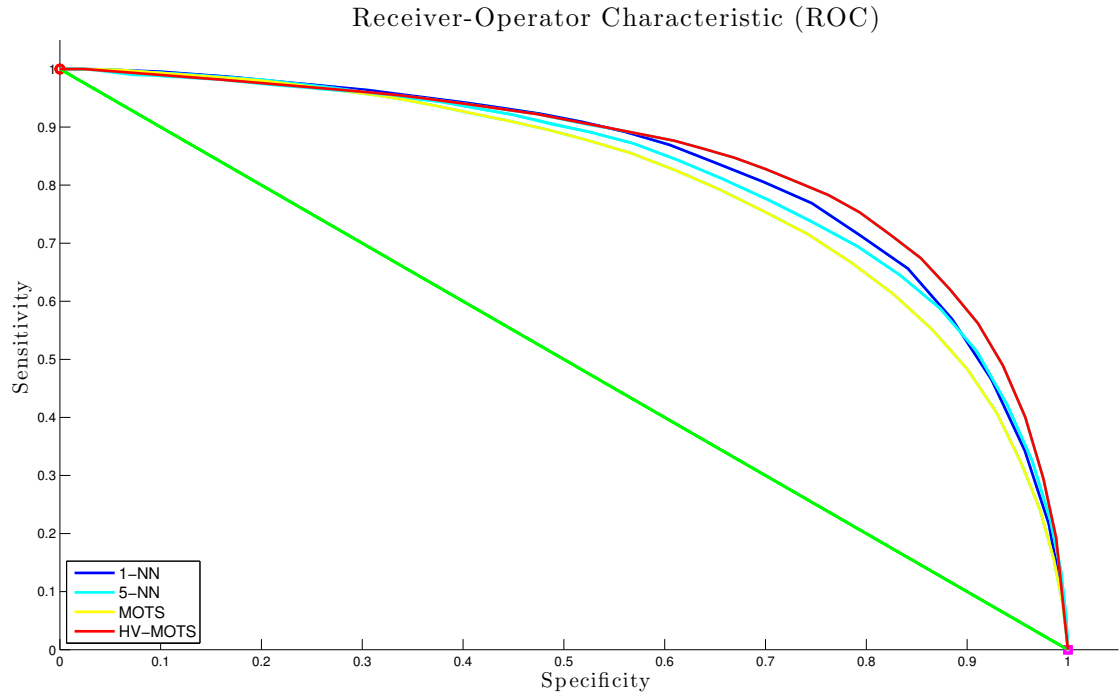


Figure 71: Receiver-Operator Characteristic (ROC) curve

	1	2	3	4	5	6	7	8	9
Rank-1	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
Rank-2	0.841	0.895	0.872	0.933	0.957	0.986	0.988	0.986	0.986
Rank-5	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
Rank-10	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013
1 st 2 nd (μ)	0.012	0.301	0.123	0.134	0.121	0.301	0.123	0.134	0.121
1 st 2 nd (σ)	0.131	0.076	0.082	0.142	0.121	0.012	0.011	0.012	0.013

Table 25: Result of *Order-3 analysis* for different levels using the HV-MOTS method

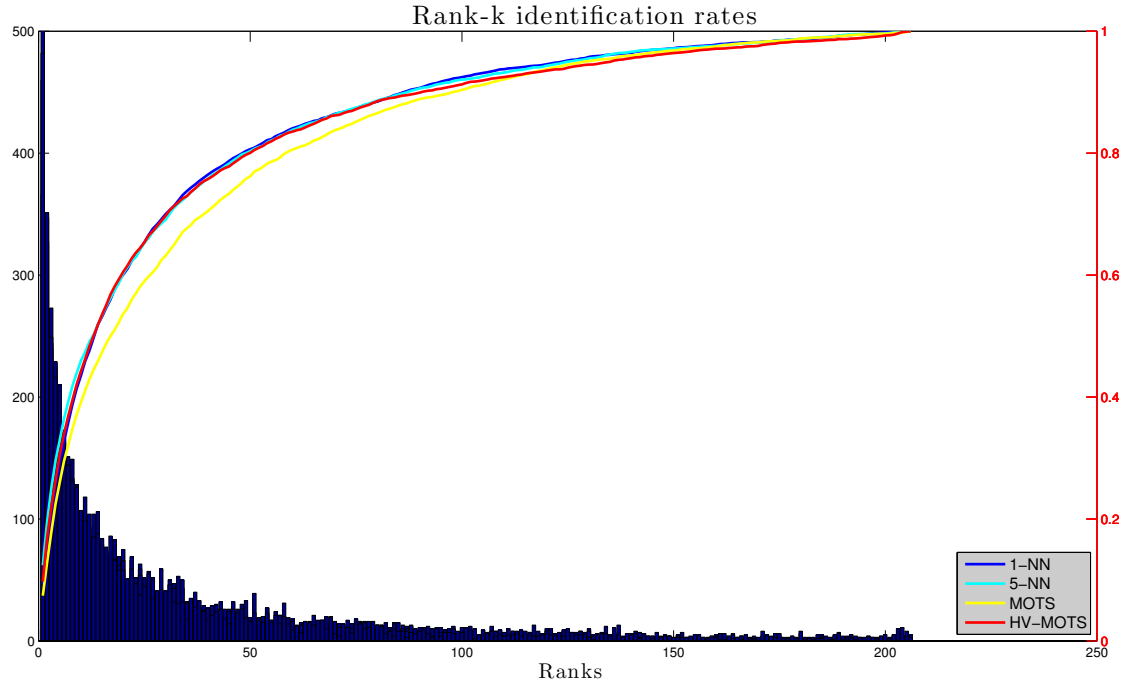


Figure 72: Results of the system identification based on the Rank-k identification rates

Menagerie analysis

The distribution of matching scores (either genuine or imposter) can exhibit *outliers* that need further investigation. As discussed by previous research [412], score distributions may vary and exhibit the presence of user groups that are fundamentally different from the distributions of the general population depending on their genuine and imposter scores. We therefore study user variations by performing a *menagerie analysis* which assess the existence of users termed as *chameleons* (simultaneous high genuine and imposter scores), *phantoms* (simultaneous low genuine and imposter scores), *doves* (high genuine and low imposter scores) and *worms* (low genuine and high imposter scores). That way, by exhibiting these singularities in users, we can further studies what makes these animals special and maybe exhibit weaknesses in the system. We will therefore look at their specific characteristic in an independent manner. Figure 73 shows the results of the menagerie analysis for the best set of features.

Comparison of datasets

We compare the results of our proposal on the two different datasets. Each of these collections has specificities that allow to study different robustness effects on our system. First, the HSCT-11 dataset allows to study the heart sounds biometric identification over large numbers of users. The baseline EER value for the HSCT-11 database is 13.66 %, obtained using the UBM/GMM method described in [38]. Our method allows to obtain an EER of XX.XX% which strongly outperforms previous research on this topic. This result confirms the feasibility of heart sounds identification. The Mars500 dataset even if limited in the number of persons has numerous interesting properties.

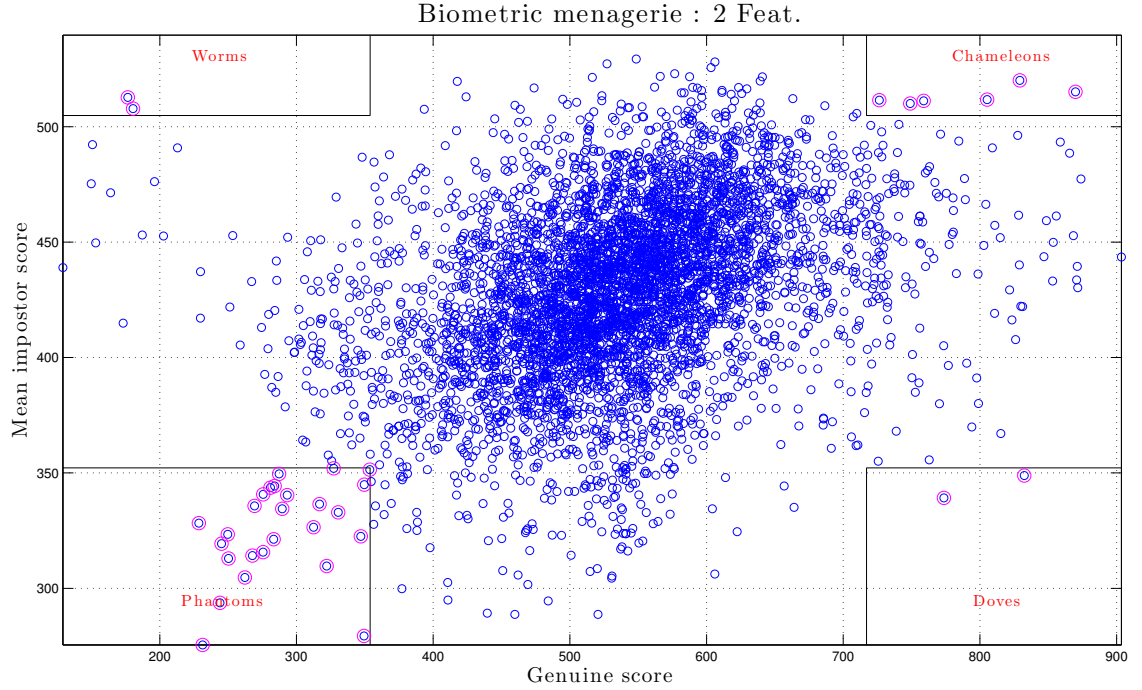


Figure 73: Results of the menagerie analysis performed over every heart beats for the best feature combination.

We obtain an EER of **XX.XX%** for this dataset. We further study the phenomenon of *template ageing* thanks to this dataset.

Template ageing

This phenomenon is known to be one of the major concern in biometrics and arguably one of the main potential weaknesses of our proposal. Indeed, the occurrence of cardiac diseases is known to strongly impact and modify the sounds produced by the heart. Furthermore, even healthy human hearts are bound to evolve with age but these long-term effects are unknown. We provide a first analysis of this phenomenon thanks to the Mars500 dataset. For each of the six subjects, recordings have been performed recurrently every three months, leading to a collection of eight samples per subjects with a maximal separation of 472 days between recordings. Our methodology is therefore to compare every possible pairs of recordings and to create *classes of ageing* for similar separations between collection. This leads to a total of seven ageing classes, each of them being three months further apart. Table 26 summarizes the distribution of these samples, number of testing instances and overall EER obtained for each class.

Parameters influence

We evaluated the influence of all the parameters implied in our heart sounds biometry system. However, for the sake of clarity, we only include here the most important parameters. The complete analysis of all the parameters is available in Annex (cf. Section A.3).

	1	2	3	4	5	6	7
Separation							
Instances							
EER							

Table 26: Results of the analysis of the *template ageing* phenomenon.

12.4.4 Comparison to existing biometrics

Table ?? presents the comparison between current results in various biometric features. We tried to gather the most recent results of major competitions in biometrics. Results presented here comes from the *Fingerprint Verification Competition* (FVC) [112], the *Iris Exchange* (IREX) [204], the *Noisy Iris Challenge Evaluation* (NICE) [310] and *National Institute of Standards and Technology* (NIST) [383]. We chose to compare our proposal to various biometrics based on the best EER value as it allows to perform a straightforward comparison of systems performances. It should be noted that for face recognition, the NIST recommends not to use the EER as evaluation measure as it makes the system work at unrealistic False Match Ratio.

Biometrics	EER (%)	Genuine	Impostors	Reference
Iris	0.057	6000	6000	[204]
Fingerprint (Standard)	0.108	27720	87990	[112]
Fingerprint (Hard)	0.687	19320	20850	[112]
Iris (Noisy)	1.31	1876	1876	[310]
Palmprint	2.569	2800	4950	[71]
Face recognition	2.83	943	88306	[55]
Signature (Online)	2.85	6374	9378	[48]
Gait analysis	6.2	50	2450	[273]
Speech	7.9	310	95790	[383]
Signature (Offline)	9.15	6527	9360	[48]
Keystroke	10.37	133	6732	[144]
ECG	10.8	73	5256	[353]
Eye movement	27.0	32	992	[180]

12.5 DISCUSSION

We provide here a discussion on the advantages and drawbacks of the heart sounds as a biometric feature. One of the inherent problem of studying heart recordings comes from their time varying nature. Commonly used biometrics like fingerprints or face recognition offers an utmost advantage in which the feature is “time-static”, meaning that the template can be extracted from a snapshot frozen in time. On the opposite a PCG signal is inherently temporal and therefore needs time to develop. Furthermore, the system requires at least a few heartbeats to obtain better performances which therefore requires more acquisition time. When comparing the enrolment template to the identification input, their differences may also be amplified by several external

Factor	Description of the effect
<i>Demographics</i>	Children can have rapid overall physiology changes whereas older people tends to have heart conditions more often.
<i>Template ageing</i>	Changes in the users biometric pattern will vary in accordance with the delay between the enrolment and identification.
<i>User behavior</i>	Can affect the recording process through unrequired movements
<i>Physiological</i>	The physiological state (<i>physical activity, stress, tension, relaxation</i>) might impact the heart sounds
<i>Environmental</i>	Disturbances in the recording process with background sounds (such as the subject speaking) or noises
<i>Sensor and hardware</i>	Differences in the sensor quality, the pressure applied and the transmission channel used to acquire the data.
<i>Heart diseases</i>	Changes produced by diseases can produce a wide disparity over time. Chronic and temporary illnesses should be studied.
<i>Spoof attacks</i>	Can hardly be performed on the system but should be considered like pre-recorded sound inputs

Table 27: Lists of both user and environmental factors that could potentially affect the performances of the heart sound identification system.

factors. We provide in Table 27 a list of these factors. We could therefore argue on our methodology that the disparity of data is too thin. Indeed, the datasets do not fully cover the problems of physiological conditions (rest, stress, anxiety, exhaustion) that may impact the morphology of the heart sounds. The topic of heart diseases has been willingly left aside, even if it could be one of the main shortcomings of our system as it can strongly disrupt the sound signature of the same person's heart. Even if we provided a first glance on the problem of template ageing, we still do not know the effect of ageing (on the scale of several decades) over the shape and strength of heart sounds. Another problem of interest emerge from the recording process itself. Indeed, several auscultation beaches can be used to obtain a recording. The heart sound signature being filtered by the surrounding body, the choice of these beaches may also influence the identification results. Finally, because of its emerging nature, the study of heart sounds distinctiveness suffers from an evident lack of data. Testing the scalability of our study would require a larger population of study.

So we could wonder why the biometric identification is working even within all these limitations. We can provide a first answer to this question by taking a closer look at the mandatory properties of biometric features [190]. First, regarding the *uniqueness* property, it is supported by medical evidence that the physiological variability (like mass distribution) provide each PCG signal with distinctive characteristics [384]. Moreover, medical research has even been devoted to reduce the variability in heart signals between individuals for diagnosis purposes [116]. In this regards, we can hypothesize that the development of the flesh of hearts is unique to each person but also the development of surroundings areas (bones, lungs, skin). This also provide a first sketch of answer on the success of the S-FC feature as the heart sounds can be related to a

source-filter paradigm, equivalent to speech (in which case the MFCC are known to be a weapon of choice). Furthermore, using the heart sounds as biometrics shows no flaw in the *universality* property as a beating heart is mandatory for any living person. Over that, recent results [4] suggests that medical biometrics is not only *permanent* over a long period, but also allows continuous identification with successive recordings input to the system. We provided a first glance on permanence in a biometric context which also show that the heart signal is consistent over a span of at least two years. Heart sounds also offers a very low *circumvention* and should be extremely robust to spoof attacks. Indeed, it seems hardly imaginable to forge or stole an heart as it undoubtedly need its possessor to produce a sound. Furthermore, it seems also hard to be concealed or hidden as burning fingerprints. The remaining properties of *collectability* and *acceptability* would however require large scale user surveys to provide a definitive answer.

13

AUDIO APPLICATIONS

The HV-MOTS classification framework was inspired by our flexible way of processing dimensions in our auditory perception. Therefore, it seems natural to now turn our attention to audio applications of our method.

13.1 INTELLIGENT SOUND SAMPLE DATABASE

13.1.1 *Content-based audio retrieval*

Content-based audio retrieval has become a popular research field, notably through the appearance of QBH introduced by Ghias et al. [143]. Most of researches devoted to this topic are based on symbolic song databases and therefore use the notion of *pitch contour* [381], which is the sequence of relative differences in pitch between successive notes. Recently there has been works that match songs directly from audio using the *melody slope* [193, 262, 435] which is the continuous equivalent of the pitch contour. The matching process must be flexible enough to allow errors in the users query and several approaches have been proposed such as HMMs [194], dynamic programming [294] or time series matching [434].

The QBE paradigm has been proposed in order to retrieve generic audio signals. QBE is based on the idea that users could find samples similar to a given example based on its spectral properties. The first QBE system was proposed by Wold et. al [404] where sounds were represented by a vector of mean, variance and autocorrelation values of spectral features. These vectors were then compared with the Euclidean distance as similarity metric. This approach known as *Bag-Of-Features* (BOF) has been extended using larger sets of features [395] or adding *relevance feedback* [312] in which the user selects its preferred results for refinement [394]. Subramanya et al. [365] used frequency coefficients from spectral decompositions and showed the superiority of DCT. They later used the multi-resolution property of the wavelet transform [364] and showed its robustness to noise. However, the selection of coefficients yields very large vectors which may be unsustainable for massive datasets. This approach was extended in [240] by using multiple statistical values over wavelet coefficients. This allows hierarchical indexing, as proposed in [242] with a pyramidal algorithm which provides an acceleration over previous approaches. Several indexing and learning schemes have also been investigated. Li [241] proposed the Nearest Feature Line (NFL) based on the idea that in feature space, lines between similar audio clips represent continuous deformations between class properties. Therefore, comparisons with queries are made with these feature lines. This method appears to yield higher classification accuracy than the Nearest-Neighbor (NN) and Nearest-Center (NC) methods. However, computing NFLs between every sound samples seems to induce a large computing and storage overhead. One of the most popular learning schemes is SVM which has been shown to outperform NN and NC [160] in classification tasks. Other machine learning techniques like Boosting [161] or GMM [173] were studied but they seem to be outperformed by the SVM-based approach.

Regarding temporal modelization, Cai et. al [65] proposed to use templates of temporal patterns for energy, harmonicity and pitch contour. Although they showed to improve accuracy, this approach seems hardly scalable because of the relative simplicity of the patterns used. More generic temporal modelization with HMMs [426] has been proposed, where comparison of HMM likelihoods with the query allows to obtain a ranked list of results. Casey [75] proposed to use the MPEG-7 feature set with an Independent Subspace Analysis (ISA) to obtain the most salient features of a sound. He further introduced a minimum entropy method [76] to train the HMM classifier which appears to outperform classical training. However, the ISA usually yields large computational overheads. The superiority of HMM cross-likelihood ratio has been shown over GMM [386] and feature histograms [171] for class-based QBE. However these studies exhibited that all the approaches are very sensitive to noise and low-quality sounds.

Another stream of generic audio querying is *Semantic Audio Retrieval* [355], which tries to learn the connection between semantic and acoustic spaces in order to perform queries on semantic concepts rather than acoustic features. The idea is to model the semantic space as a multinomial model and use a probabilistic model to connect the related acoustic properties. This approach was implemented with a mixture of probability experts in [354]. Cano et. al [69] proposed to augment this approach with a taxonomy to avoid tag confusion and polysemy. They further used a NN classifier [68] with sounds linked to concepts. Barrington et. al [28] proposed a mapping in which each dimension indicates the relative weight of its semantic concept. Casey [77] proposed to use the Passive-Aggressive Model for Image Retrieval (PAMIR) to learn the mapping between spaces. This method performs equivalently as GMM and SVM but seems to be faster. Semantic retrieval allows to perform natural language queries and circumvents the problem of manual annotation. However, it still requires a starting set of annotated sounds. This also poses severe limitations for generic sounds properties which cannot be described objectively like synthesis sounds.

13.1.2 Datasets

In order to assess the performance of our approach, we evaluate it in classification tasks using two datasets. These datasets are organized following the same database structure presented earlier (Section 7.6). First, the reference MuscleFish dataset [404] allows to compare our approach to state-of-art methods. Second, we collected a more recent and comprehensive dataset to test how our approach scales up to wider sets of data. Both datasets are available on a supporting web page ¹ so that the results of our experiments are fully reproducible.

MuscleFish

This dataset, assembled by Wold et. al [404], has been used extensively [161, 160, 241, 328, 348] in order to compare performances of different systems. It is composed of 409 sound files which are divided into 16 classes. Complete description of the dataset is presented in Table 28. Files are single-channel Sun/Next (.au) μ -law encoded audio files quantized to 8-bit with a sampling rate of 8 kHz. Loudness levels and file lengths vary over samples with the average size of a file being about 50 KBytes.

¹ <http://repmus.ircam.fr/esling/ieee-mots.html>

Musical instruments		Effects	
Altotrombone	13	Animals	9
Bells	7	Crowds	4
Cellobowed	47	Laughter	7
Oboe	32	Machines	11
Tublarbells	19	Percussion	99
Violinbowed	45	Telephone	17
Violinpizz	40	Water	7
Speech			
Female	35	Male	17
Total			409

Table 28: Description of the MuscleFish dataset used in classification tasks. 409 sounds are divided into 16 classes.

Freesound

In order to evaluate how our approach scales up to more comprehensive datasets, we collected 2193 sounds representing 54 classes from the Freesound project ², which makes this set five times larger than the MuscleFish dataset. Complete description is presented in Table 29. Files are single and double channels, WAVE and AIFF format, quantized to a minimum resolution of 16-bit with a minimum sampling rate of 44.1 kHz. Loudness levels and file lengths vary with the average size of a file being about 310 KBytes. One particularity of this dataset is that it includes a section for synthesis sounds, which are classified based on their temporal morphology.

Evaluation methodology

The goal of the classification task is to input a sound file into the system which tries to find which class it belongs to. As our method does not require any training, we use the *Leave-One-Out* evaluation methodology. That is, each file is first withdrawn from the dataset and then input for classification with the remaining set acting as a database. In order to measure performances, we use the *classification accuracy* defined as $\mathcal{A}_{cl} = N_{true}/N$ with N_{true} the number of clips correctly classified and N the total number of clips in the dataset. In order to compare different methods, we performed large-scale experiments by testing combinatorial possibilities among available descriptors. We therefore started by testing classification accuracy for every single descriptor listed in Table 2. We then tested classification with every combination of two descriptors, and so forth. Given that this testing methodology implies an exponentially growing number of tests, we keep only the top performing half of the descriptors after each step, based on their classification accuracies. We repeat this procedure and halve the set of available descriptors (in which to choose the objectives for classification) until the number of remaining descriptors is less than the number of objectives.

² <http://www.freesound.org>

Western instruments		Indian instruments		Animals		Effects	
Alto-flute	85	Santoor	15	Birds	29	Applause	41
Bassoon	80	Singing-Bowl	8	Cat	26	Footsteps	33
Cello	59	Tabla	18	Dog	21	Gunshots	15
Clarinet	76	Tambura	16	Horse	18	Heartbeats	23
Contrabass	66	Thumb-piano	17	Synthesis		Laughter	141
Glockenspiel	17	Drums		Bassline	99	Musicbox	28
Guitar (dist)	16	Crash	43	Reese	47	Paper	17
Oboe	113	Hi-hats	25	Vocoder	43	Siren	37
Saxophone	27	Kick	27	Wobble	39	Subway	11
Trombone	42	Loops	25	Speech		Sword	21
Trumpet	35	Snare	42	Female	96	Telephone	16
Tuba	74	Scratch	20	Male	87	Thunder	29
Viola	58	Toms	9	Robotic	31	Water	43
Violin	98	Tone	8	Scream	32	Whistle	26
						Zipper	17
						Total	2193

Table 29: Description of the Freesound dataset collected specifically for our study. 2193 sounds are divided into 54 classes.

13.1.3 Results analysis

We present in Table 30 the classification accuracies on the MuscleFish dataset for a growing number of objectives. For a given number of objectives, the value on the left is the *mean* accuracy over every combination and the value on the right is the *best* score obtained by a single combination. As we can see HV-MOTS consistently outperforms the other approaches in classification accuracy. This result is confirmed by the accuracies obtained on the Freesound dataset, presented in Table 31. Even with up to five times more classes and sounds, HV-MOTS exhibits an almost equivalent classification accuracy and still outperforms other methods.

More interestingly, it seems that HV-MOTS strongly outperforms other approaches in *mean* classification accuracy. This implies that given any set of features, the multi-

	1		2		3		4		5		6	
1-NN	27.2	81.7	48.3	89.2	51.4	90.9	54.7	91.7	59.1	91.7	66.3	90.5
5-NN	27.9	80.3	46.9	86.3	49.5	88.0	52.8	89.2	55.9	89.5	61.8	86.8
MOTS	27.2	81.7	51.7	85.8	56.6	83.9	60.5	82.2	61.5	80.9	62.2	75.8
HV-MOTS	27.2	81.7	55.2	91.7	68.4	93.9	77.1	94.4	83.9	95.1	89.7	95.4

Table 30: Classification results on the MuscleFish dataset for a growing number of objectives. For a given number of objectives, the left column indicates the mean classification accuracy and the right column indicates the best classification accuracy.

	1		2		3		4		5		6	
1-NN	32.3	75.8	39.4	86.1	52.2	88.6	57.8	90.6	59.4	89.3	62.4	89.3
5-NN	32.4	67.0	38.8	85.8	51.3	88.4	56.9	89.3	57.8	85.9	60.6	86.1
MOTS	32.3	75.8	41.7	83.3	54.3	82.5	59.0	81.7	61.1	71.0	57.2	67.0
HV-MOTS	32.3	75.7	47.9	89.6	65.1	90.9	82.3	92.1	87.2	93.1	89.2	93.3

Table 31: Classification results on the Freesound dataset for a growing number of objectives.

	1		2		3		4		5		6		7	
1-NN	2.4	-0.7	2.2	-6.9	2.1	-17.0	2.0	-22.4	2.3	-24.9	2.6	-23.4	2.6	-25.8
5-NN	2.7	0	1.7	-8.3	1.5	-18.9	1.4	-24.3	1.4	-28.0	1.5	-27.9	1.5	-30.3
MOTS	2.4	-0.7	2.5	-3.5	2.5	-11.8	2.5	-16.6	2.3	-22.5	1.8	-27.5	1.8	-30.9
HV-MOTS	2.4	-0.7	3.5	0	3.9	0	3.9	0	3.9	0	4.0	0	4.0	0

Table 32: Significance tests between various methods for a growing number of objectives across both datasets. For a given number of objectives, the left column indicates the mean column rank (from Tukey-Kramer HSD over Friedman’s ANOVA) and the right column gives the statistical mean difference in accuracy with the top performing method (from a one-way ANOVA).

objective approach will obtain better results. To support this claim, we provide in Table 32 the results of statistical significance tests between methods and across datasets, to rule out the effect of a particular data distribution. We use Tukey-Kramer Honestly Significant Difference (HSD) test [96] over the results of Friedman’s ANOVA to see if one method is statistically significantly different from the rest. We also present the statistical mean difference in accuracy over every combinations from a one-way ANOVA. We can see in this table that the mean column ranks and statistical mean difference in accuracy of HV-MOTS are strongly superior. The column rank corresponds here to the ranking of methods based on their accuracy results. This means that after two objectives, the HV-MOTS method is almost always in first position for any descriptor combination if ranked against other methods based on their accuracy score. Furthermore, the mean differences in accuracy increase with the number of objectives. It seems that the multi-objective classification is able to maintain the discriminative power of the best feature involved, whereas mono-objective selection will be confined by the worst features. This may go against the hypothesis that the feature set is more important than a particular learning scheme [268]. Furthermore, it seems here that the behavior of the whole class with respect to the input may be more important than the position of the input relative to the elements of the class. Therefore, even with lower dimensionality involved, the multi-objective paradigm is able to achieve a good classification accuracy. We can see that the MOTS paradigm (based on Pareto cardinality) is superior in mean classification accuracy to mono-objective selections for low dimensionality but starts to regress after four dimensions. This may come from the fact that increasing number of dimensions creates more inclusive Pareto fronts which deludes the cardinality indicator. For other methods, performances stabilize and even regress slightly after five dimensions are involved.

We present in Table 33 the confusion matrix of the best classification accuracy (95.4%) obtained by HV-MOTS on the MuscleFish dataset. The corresponding descriptor com-

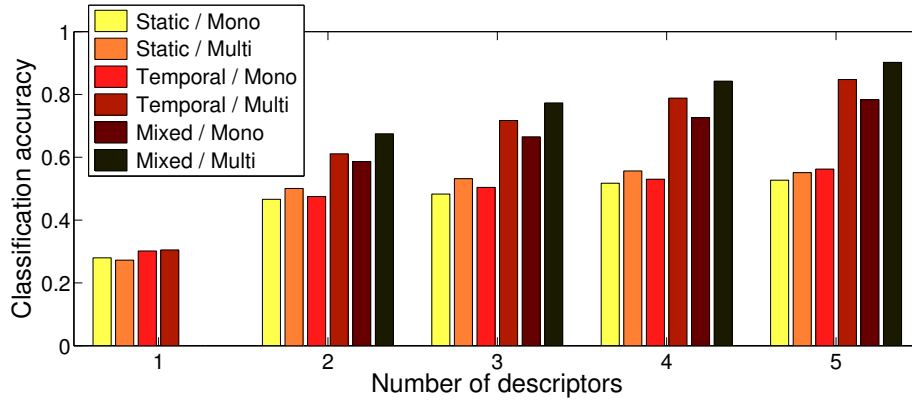


Figure 74: Comparison of the classification accuracy of using *only* temporal features, *only* static (mean and deviation) features or *mixed* sets of information with either multiobjective or mono-objective selection.

bination is composed of MFCC, MFCCDeltaStdDev, PerceptualSlope, ChromaDeltaStdDev, RelativeSpecificLoudnessDeltaStdDev and PerceptualDecrease. It is interesting to note that most of the features used are related to the temporal behavior of the sound spectrum. The descriptors which are not temporal shapes are deviations of derivative, which in fact summarize the quantity of temporal variations for these descriptors. Furthermore, this combination contains descriptors for each structural aspects of sounds, namely energy (*Loudness*), harmony (*Chroma*), spectral shape (*MFCC*) and perceptual descriptors. It should be noted that the same accuracy was obtained by 18 similar combinations (which further confirms our intuition that HV-MOTS is able to retain the discriminative power of the best features involved). If we look at the distribution of the confusion matrix, we can outline different types of errors made by the system. First, the *class similarity* errors that can be expected when similar classes are part of datasets with widely diverse class types. For instance, elements of *male speech* are confused for *female speech* and the same apply to *violinbowed* confused with *cellobowed*. Second, the *morphological similarity* errors can be observed when the spectral behavior of two classes are alike. For instance, *violinpizz* are confused with *percussions* because of the impulsive nature of such sounds. The same applies to *machines* confused with *water* because of the long-term repetitive patterns that emerge from both. Finally, in both types can be found some *reciprocal errors* where the error applies symmetrically to two classes.

The HV-MOTS method was designed based on the hypotheses that temporal shapes would improve static information and at the same time multi-objective selection would provide a perceptually more relevant and therefore more accurate classification. In order to analyze these hypotheses, we confront different views on experimental results. Figure 74 provides a comparison of the classification accuracy of using *only* temporal features, *only* static (mean and deviation) features or mixed sets of information. As we can see, the use of temporal features performs better than static features. More interestingly, it appears that best results are obtained by mixed sets of information, which indicates that normalized temporal shapes and static information are complementary sets of information. Finally, for any type of descriptors used, multiobjective selections performs consistently better than mono-objective approaches.

	Altotrombone	Animals	Bells	Cellobowed	Crowds	Laughter	Machines	Oboe	Percussion	Speech (female)	Speech (male)	Telephone	Tuba	Violinbowed	Violinpizz	Water
Altotrombone	13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Animals	0	8	0	0	1	0	0	0	0	0	0	0	0	0	0	0
Bells	0	0	6	0	0	0	0	0	0	0	0	0	0	0	1	0
Cellobowed	0	0	0	47	0	0	0	0	0	0	0	0	0	0	0	0
Crowds	0	0	0	0	4	0	0	0	0	0	0	0	0	0	0	0
Laughter	0	0	0	0	0	7	0	0	0	0	0	0	0	0	0	0
Machines	0	0	0	0	0	0	8	0	0	0	0	0	0	0	0	3
Oboe	0	0	0	0	0	0	0	32	0	0	0	0	0	0	0	0
Percussion	0	0	0	0	0	0	0	1	97	0	0	0	0	0	1	0
Speech (female)	0	0	0	0	0	0	0	0	0	35	0	0	0	0	0	0
Speech (male)	0	0	0	0	0	0	1	0	0	3	12	0	0	0	0	1
Telephone	0	0	0	0	0	0	0	0	1	0	0	16	0	0	0	0
Tublarbells	0	0	0	0	0	0	0	0	0	0	0	0	19	0	0	0
Violinbowed	0	0	0	1	0	0	0	0	0	0	0	0	0	44	0	0
Violinpizz	0	0	0	0	0	0	0	0	2	0	0	0	0	0	38	0
Water	0	0	0	0	1	0	1	0	0	1	0	0	0	0	0	4

Table 33: Confusion matrix for the best classification accuracy (95.4%) obtained by HV-MOTS on the MuscleFish dataset. The descriptor combination used is composed of MFCC, MFCCDeltaStdDev, PerceptualSlope, ChromaDeltaStdDev, RelativeSpecificLoudness-DeltaStdDev and PerceptualDecrease.

13.1.4 Comparison to state of the art

We compare our results to the state-of-art methods proposed with the same evaluation framework, namely a classification task on the MuscleFish dataset with a *Leave-One-Out* methodology. This allows to report published classification accuracies as a baseline for comparison. In their original study, Wold et al. [404] proposed to compute the mean, variance and autocorrelation of loudness, pitch, brightness and bandwidth, which together with duration amounts to a total of 13 features. By using a 1-NN rule, comparing the query to all feature vectors in the database with the Euclidean distance, they reported 80.9% classification accuracy. Guo et al. [161] later tested the applicability of a machine learning technique called *Boosting* based on a vector of 8 perceptual cepstral features which provided 78.3% accuracy. Guo and Li [160] proposed to use SVM on the same feature set and obtained 89% accuracy. Li [241] introduced the NFL method which was shown to provide 90.22% accuracy. Reyes-Gomes and Ellis [328] studied the use of GMM-EM and HMM with low entropy learning and obtained 89.9% accuracy. Finally, Shao et al. [348] used Neural Networks trained by GA with Back Propagation (BP-GA) over a set of 17 features and reported 92% classification accuracy. However their results are based on separate *Train* and *Test* sets procedure which does not allow straightforward comparison. Our HV-MOTS method allows to obtain a classification accuracy of 95.35% which outperforms previously reported accuracies for this dataset. Table ?? synthesizes the comparison between our method and previous approaches. This table also presents (right column) the minimum number of features required by HV-MOTS to obtain the same *mean* classification results. This implies that HV-MOTS is likely to obtain equivalent accuracy for *any* set of N features.

	Proposed		Equivalent	
	Accuracy	N	Accuracy	N
Guo et al. [161]	78.3 %	8	77.1 %	4
Wold et al. [404]	80.9 %	13	83.9 %	5
Guo and Li [160]	89.0 %	8	89.7 %	6
Reyes-Gomes [328]	89.9 %	-	89.7 %	6
Li [241]	90.2 %	8	91.8 %	7
HV-MOTS	95.4 %	6	-	-

13.1.5 Robustness analysis

In real-life conditions, we can expect audio collections to include sounds from different sources recorded under various conditions. Some QBE systems have been tested for robustness but usually only with regards to transcoding, using either lower sampling rates [171] or lossy data compressions [68] to simulate mobile audio databases. We test our approach by applying a wider range of distortion classes to simulate various low-quality conditions in recording

- Additive white noise resulting in 30, 20 and 10dB SNR.
- Pitch down and upconversion by 10 and 20% of pitch.
- Random signal cropping by 5, 10 and 15% of length.
- Telephone filtering with a [300, 3400]Hz bandpass filter.

	Normal	Pitch conversion				Telephone		
		-20%	-10%	+10%	+20%			
1-NN	91.69	88.02	90.71	89.73	86.06	90.95		
5-NN	89.24	85.09	87.29	87.04	83.37	88.26		
MOTS	85.82	76.53	83.37	84.60	78.97	84.11		
HV-MOTS	95.35	90.71	94.13	93.40	90.46	93.89		

	Normal	Cropping			Noise (SNR)		
		5%	10%	15%	10dB	20dB	30dB
1-NN	91.69	90.95	90.46	90.46	76.28	76.28	81.66
5-NN	89.24	88.51	88.51	88.26	74.82	74.82	78.97
MOTS	85.82	84.60	84.11	84.11	66.26	66.26	74.08
HV-MOTS	95.35	94.87	94.13	93.89	74.82	78.97	84.84

Table 34: Effects of a set of distortions on classification accuracy for different methods on the MuscleFish dataset.

These distortions are applied one at a time to each sound clip. Modified samples are then used as queries to the database (minus the original non-distorted sample) which allows comparing *Leave-One-Out* classification accuracies after distortion. We use in these tests only combinations of the best feature sets obtained in the classification task with normal quality audio. Results of the robustness analysis are synthesized in Table 34. We can see here that HV-MOTS consistently outperforms other approaches for cropping, pitch modification and telephone filtering and appears to be robust for these transformations. However, it seems that both multi-objective approaches are more brittle than mono-objective selection when considering noise robustness. Although, it should be noted that the robustness of algorithms depends on the robustness of features.

13.2 SOUND MORPHOLOGY

Our proposed HV-MOTS approach have also already been put to use in a recent work in the domain of audio perception Koliopoulou [223]. The goal of this study was to study the notion of *sound morphology* introduced by the french composer Pierre Schaeffer Schaeffer [340] in which the “*morphological profiles are meant to accurately describe the temporal evolution of some sound features and to propose an indexation and classification structure that could account for these evolutions*”. We can see that our various methods previously introduced could especially fit this framework.

13.2.1 Onset of the study

This study is inscribed in a wider framework of research which tries to tackle the evaluation of learning in interactive sound systems. The goals of this project is to learn typologies of gestures associated with environmental sounds that could be exploited in the *Sonic Interaction Design* field, which is directed at improving everyday interaction with tangible objects and interfaces. Its first step is therefore aimed at evidencing the

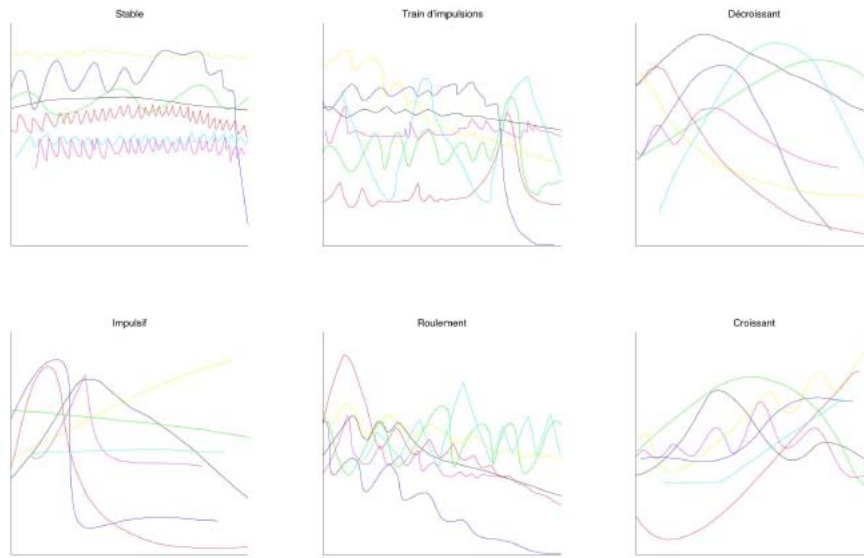


Figure 75: Temporal profiles drawn by each of the 19 participants for the set of 6 morphological classes.

morphological description of sounds by using categories of prototypical profiles for these sounds (classification) and their graphical representations (description) associated with the profiles (dynamics, harmonic) of studied signals. The latest results in this topic Misdariis et al. [272] has evidenced the pertinency of this approach in the case of dynamic profile where morphological classes have been formalized and then associated with prototypical profiles, symbolized graphically.

A first experimental study of the perception of morphological profiles was therefore realized on a set of 55 environmental sound. This study was meant to define classes of morphological profiles (dynamic and melodic) adapted to these sounds and then to conceive a formalism to describe these profiles. Finally, the study aimed at implementing a model of calculus from the temporal features in order to obtain the best signal representation for these classes. Therefore, the experiment was separated over two phases.

- Participants were asked to freely classify the sounds into 6 classes representing Schaeffer model, in order to regroup the sounds in each of the classes.
- The participants then had to draw the temporal profiles corresponding to each of the classes obtained after clustering.

The first step allowed to obtain a coherent classification between different subjects by using a hierarchical clustering approach. The second part of temporal evolution drawing for each of the participants allowed to obtain a set of morphological evolutions for each of the classes as presented in Figure 75.

The goal of the application of the HV-MOTS dataset for this problem was to evidence a possible coherence between the classification obtained by the subjective selection provided by users and the underlying audio features. The other aim of this study was to find the best set of features that could provide an objective explanation of the morphological classes.

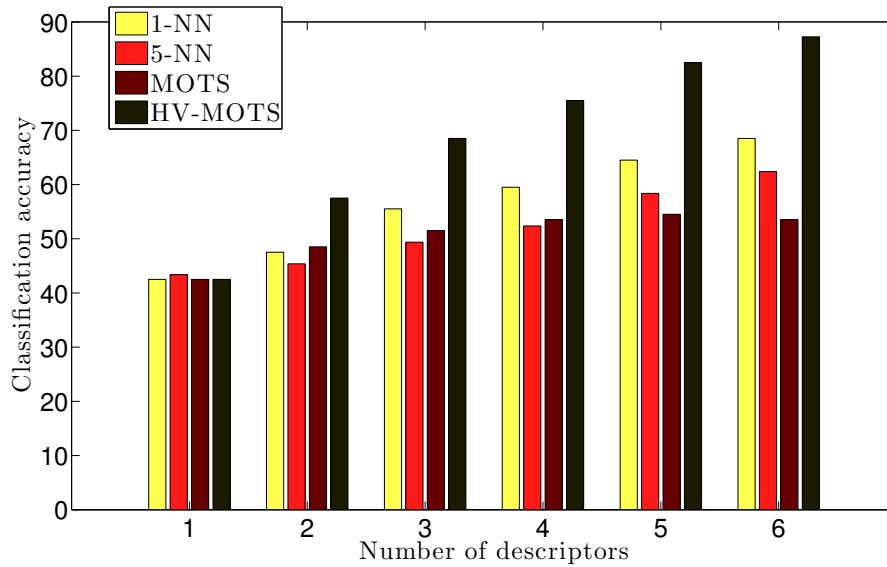


Figure 76: Accuracy of different classification methods for an increasing number of descriptors over the sound morphology.

13.2.2 Results

In order to exhibit the best set of audio features that could explain the subjective classification, the dataset of audio files has been processed by the same workflow used for analysis in the audio database application (cf. Section 13.1). Therefore each sound of the dataset is processed with IRCAMDescriptor in order to obtain all the audio features available (*Energy*, *Harmonic*, *Noise*, *Spectral* and *Perceptual*). The resulting time series are then normalized with *zero mean* and *unit variance* and resampled to a fixed length of 128 time points. Then the classification of the dataset given the subjective class is studied using the *Leave-One-Out* evaluation procedure for every combinations of audio features compared thanks to the *Dynamic Time Warping* (DTW) distance measures. The results of the classification for a growing number of objectives is presented in Figure 76.

As we can see, the HV-MOTS classification strongly outperforms all the other approaches. If we compare these results to those of the large scale study (cf. Section 11.3), it seems that the margin between 1-NN and HV-MOTS accuracies are even larger in this particular dataset. The maximal accuracy obtained at level 6 is 69.09% for 1-NN as opposed to an accuracy of 87.27% for HV-MOTS. This result correlates well with the hypothesis on which our proposal have been made. Indeed, the HV-MOTS approach was constructed to mimic our auditory perception, by producing a multidimensional organization of temporal features. Therefore, its strong superiority on an audio perception problem further confirm the validity of our hypotheses. The second goal of this study, was to obtain the set of audio features that could best explain the subjective classification and therefore provide a first glimpse on the perception of sound morphology. Table 35 presents the best features combination obtained by the classification step.

As we can see, the best combination provides a set of audio features that embed several aspects of the sound properties. Two descriptors describe the *energy* components of sound (*energy envelope* and *loudness delta*). The *fundamental frequency* feature allows to describe the evolution of *pitch*. For the *harmonic* content, the inharmonicity describes

Best combination obtained (Accuracy : 87.27%)	
<i>Energy Envelope</i>	42.51 %
<i>Spectral Rolloff</i>	37.82 %
<i>Loudness Delta</i>	28.83 %
<i>Noisiness</i>	21.13 %
<i>Inharmonicity</i>	19.24 %
<i>Fundamental Frequency</i>	18.65 %

Table 35: The best combination (6 features) obtained thanks to the HV-MOTS classification paradigm provides a classification accuracy of 87.3%. The right column shows the individual classification accuracy for each feature.

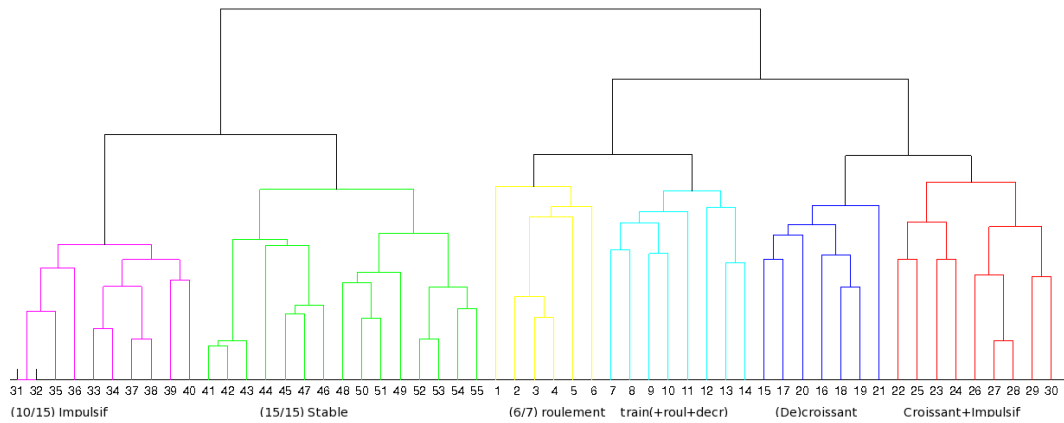


Figure 77: Results of a hierarchical clustering performed on the sounds of the study using the audio features selected thanks to the HV-MOTS classification paradigm.

the evolution of harmonic peaks while the *noisiness* gives the ratio between the noise and harmonic components of a sound. Finally the *spectral rolloff* provides a description of the overall position of the energy in the spectrum. To validate this set of features, Figure 77 presents the results of the hierarchical clustering performed on the sounds of the study using the audio features selected thanks to the HV-MOTS classification paradigm. This clustering have been performed using the distance between sounds given the set of audio features. The clustering is performed by computing the cityblock (L_1) distance between pairs of features and then the shortest distance between pairs is used to regroup sounds in clusters.

The selected set of features allows to obtain almost the same clustering as obtained from the subjective ratings, which further confirms its validity. As we can see in this figure, three of the classes are perfectly clustered. For the other classes, almost all sounds have been clustered in their corresponding set. This study exhibit the validity of the HV-MOTS paradigm in an audio perception framework. Our approach allowed to obtain a strongly higher classification with a set of features which is then validated through studies. The remaining part of the study in Koliopoulou [223] exhibits that this selected set allows to compute a meaningful set of prototypical profiles for each

morphological class. We do not include this part for the sake of clarity and redirects the interested readers to Koliopoulou [223].

CONCLUSIONS

Part V

GOING BACK TO MUSIC

15

ORCHESTRATION

We now get back to our original artistic problematic. As we discussed in Chapter 1, musical orchestration is an art of musical writing that rely on the spectral characteristics unique to each instrument. Therefore, it goes beyond the realm of symbolic processing and thrives in the territory of signal possibilities. Trying to tackle the problem of orchestration from a scientific angle involves the use of unformalized knowledge and unveils numerous facets of complexity.

15.1 ON THE COMPLEXITY OF ORCHESTRATION

The history of music is littered with critical moments when the established conventions are no longer sufficient to accompany the novel musical trends. Therefore, when scientific research focuses on musical issues, there is no choice but to follow this need for complexity. We will try to highlight some open problems in different areas of complexity in the orchestration topic. We divide this discussion into three major areas of complexity: *combinatorial* complexity, *temporal* complexity and *timbre mixtures* complexity.

15.1.1 *Combinatorial complexity*

The first facet of orchestral complexity arise from the exponential number of possible instrumental mixtures. Indeed an orchestra can count up to hundreds of players and is usually composed of a large-scale assortment of instrumental groups. Each instrument is able to produce a variety of playing modes on a vast range of notes which can be played at various intensities. Therefore, trying to find every combination of timbre that can be played by an orchestra implies to solve a NP-Complete problem. Furthermore, it is hard to predict the properties of every instrumental mixture, as it is computationally intense to compute the signal features on each mixture. We can take a quick sobering experiment to illustrate this problem. If we suppose that each instrument can play only 25 notes (a grand piano can play up to 96 notes) at 2 dynamics (the symbolic notation contains 9 dynamic symbols but crescendos can produce continuous variations) with only 2 playing modes. This yields 10^2 possible musical atoms for each instrument. Hence, a set of N instruments would lead to 10^{2N} possibilities. Now, if we suppose that our computer is able to evaluate the spectral features of a combination in only 10^{-9} s (1ns is extremely fast considering that this computation is currently at most real-time), then evaluating all combinations requires

$$N = 5 \quad 10^{10} \text{ combinations} \rightarrow 10 \text{ seconds}$$

$$N = 8 \quad 10^{16} \text{ combinations} \rightarrow 115 \text{ days}$$

$$N = 10 \quad 10^{20} \text{ combinations} \rightarrow 3.170.979.198.376 \text{ years}$$

This small experiment allows to get a glimpse on the extent of this problem. Furthermore, we have taken here extremely simplifying assumptions. It is therefore mandatory to find an approach that can handle this combinatorial explosion.

15.1.2 Temporal complexity

As we discussed in Section 2.4, the temporal structures are of prime importance in music processing. However, if we look back at Figure 1, reaching back to the scale of orchestration reveals a new dimension of complexity. Indeed, the *macro-temporal* evolutions combine vertical writing (the arrangement of different voices in a restricted time frame) and horizontal writing (the development of the complete musical structures over time). Therefore, orchestration unveils an interaction between the *micro-temporal* properties of musical atoms and the *macro-temporal* articulations of their organization. Each scale relates to different levels of complexity. We have studied the micro-level extensively throughout this document. However, the macro-level is a necessary condition for the orchestration of polyphonic sequences. As we will discuss in Section 15.2.2, we address the vertical and horizontal dimensions of orchestration simultaneously. Therefore, we must face the problem of macro-temporal timbre which evolves continuously. It seems impossible to consider this dynamic orchestration as a mere extension of a static model in where the notion of time would come down to a series of segmented instants. Hopefully, this is where all the knowledge gained in previous chapters will prove its usefulness.

15.1.3 Orchestral timbre

The advent of computer music has allowed to explore a novel universe of sound possibilities. The electronic instruments, introduced "unbelievable" timbre, profiling a new dimension in the imaginary of sound. These are now an integral part of contemporary musical discourse. The timbre gradually became a central element of this new musical language and is now even looked upon as the backbone of musical writing. The large number of instruments and their variety of timbre offer a virtually infinite palette of orchestral colors. However, as we discussed in Chapter 1, there are numerous problems when trying to cope with the multidimensional aspect of timbre and the prediction of sound mixtures properties. One of the main facet of complexity in this context is to understand the coupling between instruments that arise from sound mixtures.

THE PHENOMENA OF EMERGENCE These phenomena are still poorly understood today but are embedded in the study of orchestration. Indeed, it is well-established that two instruments create a timbre different from the simple sum of their respective timbre. It is the relationship between instruments that create different timbres, because of their complex acoustic interactions. Therefore, these phenomena can cause the appearance of elements in the spectrum of a mixture that are not in any of its constituents. Reciprocally, the *masking* and *phase* effects can cause the disappearance of components because of their relative levels or closeness of spectrum (which illustrates the idea that the mixture spectrum does not follow a linear addition). The phase effect has been extensively studied and also seems deeply involved in recognition of the instruments. Other interesting phenomena can emerge from an accumulation of sound sources such as the phenomena of *unison* and *chorus*. Conversely, the instrumental *fusion* can create the illusion that two instruments merge their sounds to the point that they are impossible to discriminate.

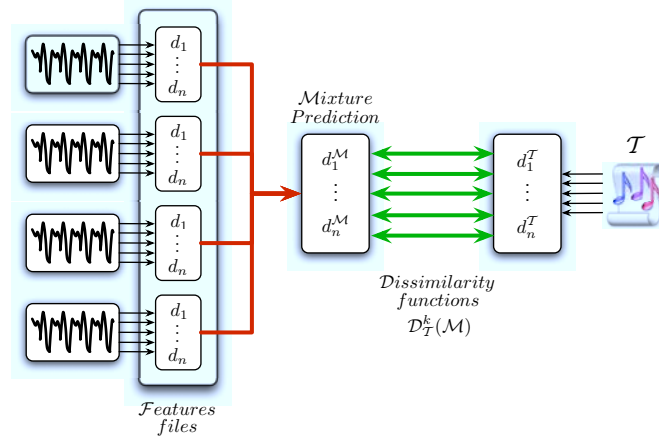


Figure 78: Original approaches to tackle the problem of computer-aided orchestration. A target sound file is fed to the system which will try to reconstruct it by using a set of instruments. The system rely on a set of feature files which are combined through prediction functions.

15.2 HOW TO REIFY ORCHESTRATION

After seeing these facets of complexity, we could wonder how to approach the orchestration on a scientific basis. Automatic orchestration is a very recent topic in computer-aided composition, for which only few systems exist today. In order to find a single point of entry to this problem, the systems all rely on the concept of a *target* to be optimized. An orchestration problem then consists in finding a mixture of instrument samples that minimizes different spectral distances. This approach is illustrated in Figure 78, where the goal is to find the mixtures of instruments that best match a given timbre. The user specifies the instruments he would like to use (*constraints*) and the characteristics of the sound to be produced (a sound *target*). This target is a sound file which defines the audio features to reproduce. Then, an orchestration engine relies on an instrumental *features database* to suggest instrument combinations (orchestration proposals) that sound close to the target. The problem is therefore to find an efficient way to converge towards closely sounding elements (and hence circumvent the combinatorial complexity).

15.2.1 Existing systems

Psenicka [311] first developed a system called SPORCH (SPectral ORCHestration) where the search is based on an iterative matching pursuit algorithm. Each instrument in the database is associated with a series of pitch, dynamics and a collection of the perceptually most significant harmonics. The target is analyzed with a spectral analysis procedur. The algorithm first selects the best matching instrumental sample. The spectrum of this sample is then subtracted from the target and removed from the database. The algorithm iterates while trying to minimize the Euclidean distance between the target and the current mixture. Rose and Hetrick [332] later proposed a tool to analyze existing orchestrations and propose new ones. The instrumental knowledge is summarized by the average harmonic spectrum calculated on the sustained part of the instruments. The algorithm is based on a Singular Value Decomposition (SVD)

of the spectrum contributions. The target is finally expressed as a weighted sum of the spectra contained in the database. Another system was subsequently proposed by Hummel [183]. An iterative algorithm is also used, but the distance is based on the spectral envelope rather than the harmonic partials. The system calculates the spectral envelope of the target and repeatedly search for the best approximation. Consequently, the harmonic structures can be very different between the target and the resulting mixture. The author advise to use his system with sounds lacking pitch such as whispered vowels.

Recently, an interesting approach was proposed by Carpentier [72]. Although the formulation of the problem remains the same [73], the system is based on multiobjective genetic exploration [375]. This allows to propose a set of optimal solutions, rather than a single instrumental combination. The target is represented by a set of audio features that describe different aspects of the sound to be reproduced with a pre-determined orchestra. Using a genetic algorithm, the system retrieves the Pareto front of solutions. The user then selects a solution, which allows the system to infer its preferences among the different dimensions. The instrumental knowledge is based on the model developed by Tardieu [374]. This model uses Gaussian Mixture Models (GMM) to learn the distribution of features for a large number of samples and allows to infer the properties of missing ones. Furthermore, a set of functions [373] were developed for each feature which can predict the characteristics of a mixture of instrumental sounds.

15.2.2 Discussion

We try to highlight here the shortcomings of existing orchestration systems.

- The main goal of every system is to approximate a *sound target*. This restricts their scope of study to a very narrow case of orchestration. Indeed, most of the time, composers does not have a well-formed example of the sound to obtain. The orchestration generally aims to produce, not to re-produce a timbre.
- These systems make the implicit assumption of linear additivity of the timbres. However, the simple consideration of the phase effect is enough to put it into default. The predictive capacity of an additive model may be doubted, despite its computational advantages.
- Finally, the most important shortcoming of every previous system is that they only provide vertical orchestration by focusing the analysis on sustained instruments. Therefore, even the best proposed solutions completely neglect the temporal evolution of the sound target. The instrumental knowledge is, therefore, limited to sustained harmonic sounds without any temporal variations.

As all previous systems rely on “time-blind” features, they are only able to provide static orchestrations. However, as we discussed in Section 2.4 the territory of timbre is not confined to a static structure of proportions. It rather comprises “variation laws” in a context continuously evolving over time. It is therefore essential to move to a higher level of modeling, by assessing the complete temporal structure of spectral features. Advantages of this approach are twofold. First, the generated orchestrations are improved and more realistic by reproducing the whole spectro-temporal structure. Second, it allows the use of evolving playing modes like crescendo, glissando, multi-phonics and so on. For all these reasons, we focus on providing a system that addresses the problem of time in musical orchestration. However, we will also tackle the two

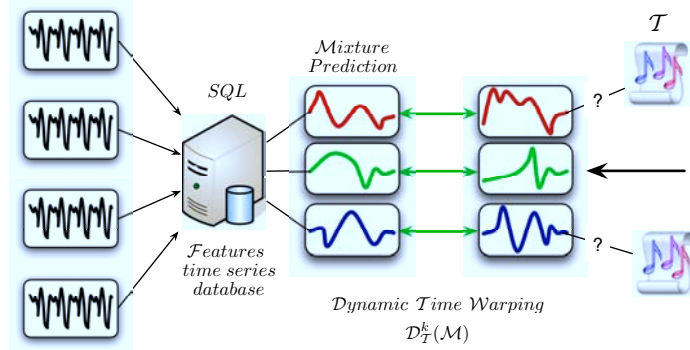


Figure 79: A new approach for computer-aided orchestration. Unlike previous systems (cf. Figure 78), the comparison between targets and sound mixture is based on their temporal evolutions. The knowledge is encapsulated in a SQL database in order to provide an almost infinite source of knowledge. Finally, the need of a well-formed audio example is bypassed by the direct input of temporal shapes.

other aforementioned shortcomings. Hence, we introduce the notion of *abstract target* and also bypass the additive model by using *prediction functions*.

15.3 GOING FURTHER IN COMPUTER-AIDED ORCHESTRATION

We outline here our main goals when devising an orchestration system that could take into account the temporal evolution of audio features. Unlike previous systems (cf. Figure 78), the comparison between the target and sound mixtures is based on the *temporal evolution* of their audio features. The knowledge is encapsulated inside a SQL database in order to provide an almost *infinite source of knowledge* (without loading all features in live memory). Finally, the need of a well-formed audio example could be bypassed by providing *abstract targets* such as the direct input of the temporal shapes. The final workflow embedding these improvements is shown in Figure 79.

15.3.1 Algorithmic choices

Knowledge database

The previous systems all rely on a set of files that contains the static audio features of instrumental sounds. This source of knowledge is inherently limited as it requires to load the complete database in live memory (in order to avoid the cost of disk accesses along the search algorithms). Therefore, the systems must limit the quantity of their own knowledge. This memory-based approach seems to be wildly contradictory with the combinatorial complexity that we exhibited in Section 15.1.1. Furthermore, as we intend to focus on temporal shapes, the dimensionnality of data to be processed will be greatly extended. Hence, we decided to use a SQL database architecture in order to store the audio features. The advantages of this choice are two-fold. First, it allows a dynamic and almost infinite source of knowledge. Second, by using the work previously presented (cf. Section 13.1), this database provides the MOTS, MOSEQ and QVI embedded inside the orchestral knowledge. We will see that these improvements can greatly enhance the orchestration search algorithm. (cf. Section 15.4.2).

Temporal matching

This is where all the knowledge gained from the past chapters can be put to use to improve our original artistic problematic. First, the use of the MOTS framework (cf. Chapter 7) in the database allows to perform temporal queries on the instrumental knowledge as well as using the MOSEQ and QVI paradigms (cf. Section 7.5). Then, when trying to optimize the sound mixtures to match a given target, we assess their temporal similarity with the *Dynamic Time Warping* (DTW) distance measure. Hence for each sound mixture, we forecast its audio features by using the prediction functions provided by Tardieu and Rodet [374]. Even if these functions were devised for static values, we consider that they can be extended to time series on a point-by-point basis. Based on these predicted time series, we consider the DTW as a dissimilarity measure to select the best mixtures.

Abstract target

The concept of a *sound* target restricts the scope of study to a very narrow case where composers have a well-formed example of sound to obtain. However, since we provide a new approach that can handle temporal matching, we can bypass this simplification. As presented in Figure 79, we provide several forms of targets ranging from sound files to purely abstract targets. First, the user can still input a single sound file to the system. The analysis module then show the temporal shapes of the corresponding sound, which the user is free to modify. Then, the system offers the possibility to use *multiple* sound files, by selecting the temporal shape of a first target to be combined with the feature of other sound files. The user can input temporal shapes independently of the sound files. Finally, a *purely abstract* target can be created in which all the features to optimize are the result of a direct input procedure.

15.4 ABSTRACT TEMPORAL ORCHESTRATION - MODULAR STRUCTURE (ATO-MS)

We introduce in this section a novel system for orchestration generation, called *ATO-MS*, that address the shortcomings of previous systems. *ATO-MS* is a multi-objective genetic optimization system which allows to find relevant sound combinations that match the temporal evolution of several audio features. Using multiobjective genetic algorithms allows to evolve *populations* of sound mixtures that jointly minimize a set of objective functions. These functions are defined as the distance between selected audio features. Each distance can either represent the DTW between time series of audio features, or the Euclidean distance between average features. Each *individual* in the population represents a potential sound mixture solution. The individuals contains the symbolic properties of the sound mixtures, encoded inside a *genome*. Therefore, each genome defines an orchestral score which will be used to approximate a given target. We define the orchestral genome in the following way

(Instrument ₁)			...	(Instrument _N)		
Sample ₁ ^{id}	Onset ₁	Duration ₁	...	Sample _N ^{id}	Onset _N	Duration _N
[0 ... DB]	[0 ... T]	[0 ... T - O ₁]	...	[0 ... DB]	[0 ... T]	[0 ... T - O _N]

The set of instruments to be used by the system is pre-defined by the user before starting the algorithm. Then for each instrument, the Sample_i^{id} defines the index of the

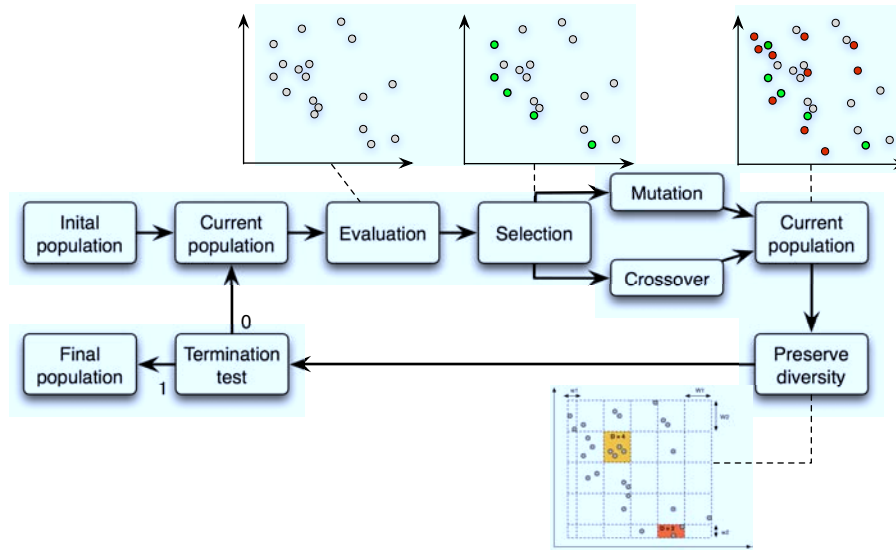


Figure 80

sample to use in the database. This sample will subsequently define the note, playing mode and intensity of the corresponding instrument. Onset_i defines the starting time of the i^{th} instrument and Duration_i defines its playing length. The goal of the genetic algorithms is then to make these individuals evolve. This is performed by trying to mimic the process of natural selection, as illustrated in Figure

15.4.1 Entropic segmentation procedure

We propose in this section a novel segmentation procedure that could account for various levels of information. In our context, the biggest flaw of segmentation procedure is that they usually only seek to define the onsets of different events (or as we termed *musical atoms*). However, in the context of orchestration, we need to “break down” further this musical unit as we seek to reconstruct it with a mixture of instruments. Therefore we need a multi-level segmentation procedure that can give access to parts of the time series that exhibit a sense of coherence. For that purpose we use the measure of *approximate (Shannon) entropy* which gives us a measure of the average information content in a segment.

Our segmentation method first relies on a *bottom-up* Piecewise Linear Approximation (PLA) (cf. Section 5.4.2). Hence, we start by considering that each pair of points of the time series as a potential event. We then work by agglomerating different segments together, starting from single points, until the complete series is considered as a single segment.

We present in Algorithm 15.1 the final implementation of the *multi-level entropic segmentation* procedure.

15.4.2 Optimal warping

We now introduce our new algorithm for computer-aided orchestration. Overall this proposal can be classified as a multi-objective genetic optimization algorithm. However, because of the central use of time series in our system, we enhance this approach by

Algorithm 15.1 The *multi-level entropic segmentation* procedure

```
multiLevelSegmentation(data, thresh, minSegs)
// Perform Piecewise Linear Approximation
data = (data -  $\mu_d$ ) /  $\sigma_d$ 
reconstruction = data
segments = 1:(size(data,1) - 1)
// Initialize the cost of merging each pair of segments
for i  $\in$  [1 ...  $N_{seg} - 1$ ]
    cost(i) = compute_error(data, merge(segments(i), segments(i + 1)))
end
// Keep merging until reconstruction error or number of segments is attained
while errRecon < thresh ||  $N_{seg} > minSegs$ 
    id = min(cost)
    segments(id) = merge(segments(id), segments(id + 1))
    remove(segments(id + 1))
    cost(id - 1) = compute_error(data, merge(segments(id - 1), segments(id)))
    cost(id) = compute_error(data, merge(segments(id), segments(id + 1)))
    errRecon =  $\mathcal{D}_{\mathcal{L}_2}$ (data, segments)
end
// Compute the entropy of ea
// Start the multi-level entropic grouping
while  $N_{seg} > 1$ 
% Compute entropy for each base segment
for i = 1:length(segments)
    segVal = tmpSeries(segments(i).lx:segments(i).rx);
    segments(i).entropy = approximateEntropy(1, 0.2 * std(tmpSeries),
    segVal);
    segments(i).deltaEnt = inf;
end
tmpQueue = segments;
hierarchicalSegment = segments;
curSegment = segments(1);
while ~(curSegment.lx == 1 && curSegment.rx == length(tmpSeries))
    deltaEntropy = zeros(length(tmpQueue) - 1, 1);
    for i = 1:(length(tmpQueue) - 1)
        if isinf(tmpQueue(i).deltaEnt)
            segVal = tmpSeries(tmpQueue(i).lx:tmpQueue(i + 1).rx);
            tmpEntropy = approximateEntropy(1, 0.2 * std(tmpSeries), segVal);
            tmpQueue(i).deltaEnt = tmpEntropy - (tmpQueue(i).entropy + tmpQueue(i +
            1).entropy);
        end
    end
    deltaEntropy(i) = tmpQueue(i).deltaEnt;
end
[v eID] = max(deltaEntropy);
curSegment = struct;
curSegment.lx = tmpQueue(eID).lx;
curSegment.rx = tmpQueue(eID + 1).rx;
```

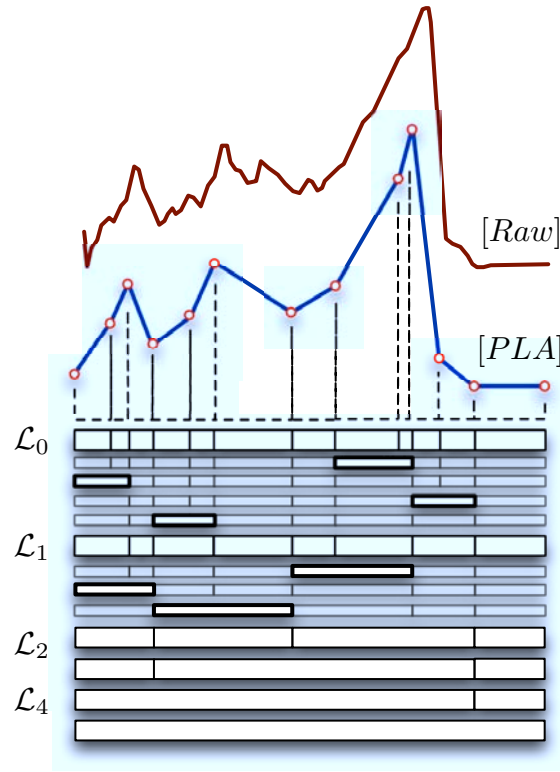


Figure 81

using a specifically designed *local optimization* method. This enhancement strongly rely on the MOTS paradigm and the entropic procedure in order to refine the proposed solutions after each iteration of the algorithm.

LOCAL MOTS OPTIMIZATION As the introduction of temporal structures implies a great increase of the problem complexity, we had to devise a way to enhance the genetic search. One of the biggest flaws of such algorithms is its relative blindness in the first population generation. Indeed, individuals are randomly chosen at the beginning to start the search procedure. Therefore we use our novel MultiObjective Time Series (MOTS) matching algorithm as a "kick-start" for the multiobjective genetic algorithm. This idea is depicted in figure 16. That way, the blindness of the initial population is reduced, as several "self-solving" efficient instruments are added prior to the search.

ALGORITHM We now introduce our complete *Optimal Warping* algorithm for solving the problem of temporal orchestration, which is detailed in Algorithm 15.2.

15.4.3 Comparison with Orchidee

In order to evaluate the proposed improvements on an algorithm for computer-aided orchestration, we will compare it to the previous best performing algorithm to serve as a baseline. In order to perform a fair and objective evaluation, we will therefore use exactly the same evaluation procedure proposed in Carpentier [72]. The various algorithms are thus evaluated on 500 monophonic problems with cardinality and

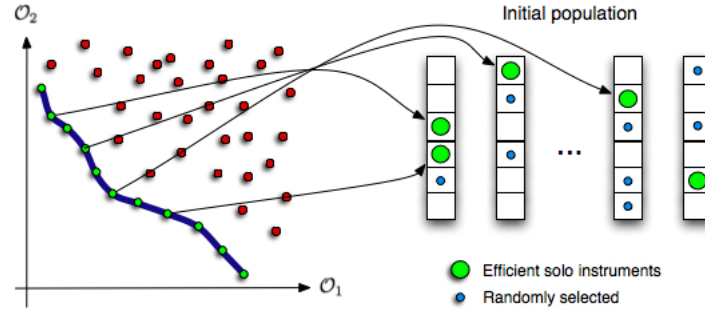


Figure 82: We combine paradigms by using the MOTS algorithm in order to find efficient solo instruments that can be used as “seeds” to the initial population. That way we reduce the blindness of pure randomness

Algorithm 15.2 Our proposed *Optimal Warping* algorithm which combines a multi-objective genetic algorithm and a local optimization procedure, based on a combination of the MOTS framework and the entropic segmentation method.

optimalWarping(\mathcal{T})

segments \leftarrow multiLevelSegmentation(\mathcal{F}_{time}^k)

seeds \leftarrow MOTS_{match}(segments)

pop \leftarrow randomPopulation(N_{pop}^{init} , seeds, onsets)

criteria \leftarrow evaluatePopulation(pop, \mathcal{T})

pareto \leftarrow extractPareto(criteria, pop)

while iter < iter_{max}

 C évaluer_{population}(X; T)

 tirer un jeu de poids

 f kCk

 f pénaliser_{fitness}(f; z)

 pool \leftarrow selectTournament(pop, N_{pool})

 pool \leftarrow crossover(pool)

 pool \leftarrow mutation(pool)

 crit_{pool} \leftarrow evaluatePopulation(pool)

 pop \leftarrow pop \cup pool

 pareto \leftarrow updatePareto(criteria_{pool}, pareto)

if N_{pareto} > N_{pareto}^{max}
 P decrease_{PADE}(P; N_{pareto} max ; random)

if N_{pop} > N_{pop}^{max}
 X X n P

 C évaluer_{population}(X; T)

 X decrease_{PADE}(X; N_{pop} max#P; f) // (. . ensuite seulement ceux de densité

 X X [P

end

end

return pareto

	Monophonic constrained		Monophonic unconstrained	
	Orchidée	Optimal W.	Orchidée	Optimal W.
Superiority	1.20 %	24.40 %	0.80 %	15.60 %
Dominance	21.60 %	98.40 %	4.80 %	95.20 %
Converge	43.60 %	56.40 %	41.20 %	59.80 %
Diversity	61.60 %	38.40 %	58.40 %	41.60 %

Table 36: Comparison of algorithms on the monophonic orchestration problems

	Polyphonic constrained		Polyphonic unconstrained	
	Orchidée	Optimal W.	Orchidée	Optimal W.
Superiority	2.40 %	14.20 %	0.20 %	19.20 %
Dominance	3.80 %	96.20 %	2.40 %	97.60 %
Converge	39.40 %	63.60 %	41.00 %	59.00 %
Diversity	61.0 %	39.00 %	57.60 %	42.40 %

Table 37: Comparison of algorithms on the polyphonic orchestration problems

pitch constraints for which the search space is known, 500 polyphonic problems with same constraints and same search space, 500 unconstrained monophonic problems for which the search space is unknown and 500 unconstrained polyphonic problems with same search space. All the algorithms use the exact same set of parameters which are $N_{pop}^{init} = 200$, $N_{pop}^{max} = 500$, $N_{mate} = 50$, $N_{pareto}^{max} = 200$, $Iter_{max} = 20000$.

In order to compare the performance of all the algorithms, we also use the same evaluation measure as proposed in Carpentier [72]. The four following measures are considered, where compared algorithms are named \mathcal{A} and \mathcal{B}

Superiority We compute the dominated hypervolume of each algorithm which is not by the other, therefore $\mathcal{H}(\mathcal{A}, \mathcal{B})$ and $\mathcal{H}(\mathcal{B}, \mathcal{A})$.

Dominance The dominance measure is a combination of the *epsilon indicator* (which factor is required in order for \mathcal{A} to dominate \mathcal{B}), the *coverage* (percentage of elements of \mathcal{B} dominated by \mathcal{A}) and the *binary hypervolume* (hypervolume dominated by each algorithm).

Converge Defined as the mean of euclidean distances between each solution of the algorithm and the ideal point.

Diversity This measure is computed as the spread of the distribution over the optimization space.

We therefore compared the original algorithm (*Orchidée*) Carpentier [72] and our *Optimal Warping* algorithm using these measures and report the results in Table 36 for monophonic problems and Table 37 for polyphonic problems.

15.4.4 Modular structure

The new system has been architected around an *Object-Oriented Programming* conceptualization. Therefore, the system in itself is an extensible and modular structure which

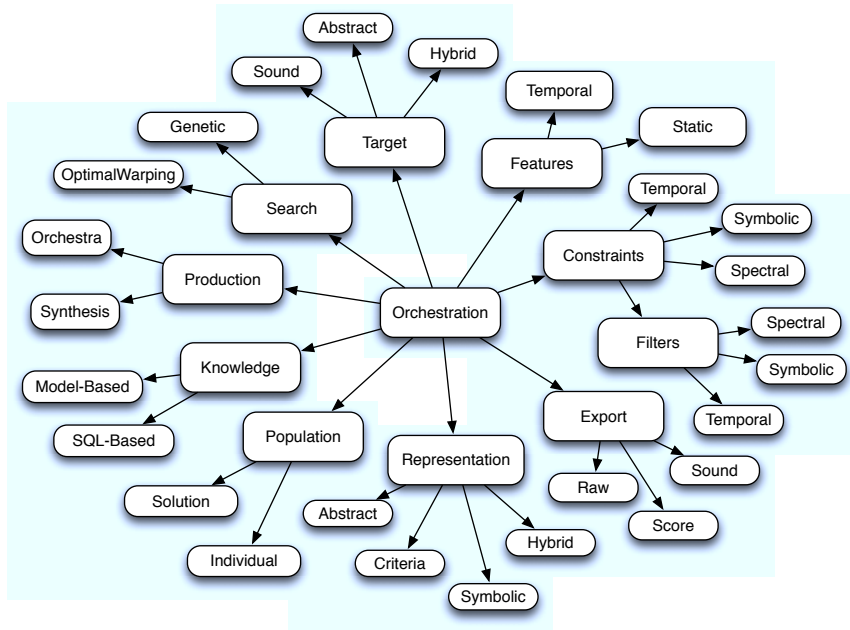


Figure 83: The current prototype for Abstract Temporal Orchestration with Modular Structure (ATO-MS) features an extensible architecture of modules to tackle the problem of Computer-Aided Orchestration.

can allow modification, extension, removal and even complete introduction of tasks and resolution methods. Figure 83 summarizes the modular structure and components of the new ATO-MS system. The *Session* object centralizes every information about a current orchestration problem and contains current instances of every sub-part of the problem. The *Production* object informs the system on the current means of sound generation that can be used, through the capacities and set of instruments that are allowed for a specific orchestration. The *Knowledge* object is used by the system to retrieve the symbolic and spectral feature for individuals used in the search process (the current SQL database system). The *Search* modules is the core of the orchestration system. These objects represent algorithms that are able to provide fast and orderly solutions, which can currently be set as a *SearchGenetic* (multi-objective genetic algorithm with an hypervolume domination criterion) or the *SearchOptimalWarping* presented in the previous sections. The *Target* defines the objectives to attain and therefore which values to approximate for each feature. As explained previously, the target can either be a *TargetSound* or a *TargetAbstract* that allow to directly input a set of features. The *Features* objects allow to define for each spectral feature its computation method, its prediction function and how to compute the distance between the same feature of two different mixture. Finally, the *Population* defines the current set of proposals computed by an algorithm to a particular orchestration problem, therefore composed of a set of *Solutions* each of which is a mixture of several *Individuals*.

15.4.5 Database

The database in our orchestration problem is a representation of the extent of knowledge that we might possess over the acoustic capacities of orchestral instruments.

15.4.6 *Interface*

Family	<i>Orchidée</i>		<i>ATO-MS</i>	
	Modes	Samples	Modes	Samples
Bassoons	8	458	18	1.359
Clarinets	18	814	37	3.959
Flutes	20	1.191	47	2.946
Horns	11	616	26	1.507
Oboes	11	695	33	2.950
Saxophones	2	161	22	517
Contrabasses	21	2.498	37	3.754
Viola	25	3.017	50	5.176
Violins	26	2.911	51	5.912
Violoncellos	25	2.930	49	4.549
Trombones	20	1.187	32	3.481
Trumpets	18	795	41	2.191
Tubas	9	891	21	2.064
Total	215	18.164	464	40.365

Table 38: Comparison of the orchestral databases used as a knowledge source for *Orchidée* (left) and our *ATO-MS* system (right).

16

OTHER ARTISTIC APPLICATIONS

16.1 MOSEQ INTERFACE

16.2 QVI INTERFACE

16.3 IPAD INTERFACE

All the aforementioned technical characteristics have been implemented as a multitouch interface using the iPad. This interface embed all interaction schemes and capabilities in an OpenGL / Objective-C framework. It has been developed as a client to a local computer server (which allows an extensive storage size) using OSC communication but can also be self-contained (even if the database size is inherently limited). The main screen of the interface is showed in figure 6.

16.4 SPECTRAL MAQUETTES

We present here a preliminary system that could embed several composition paradigms. Its goals is to allows a link between the micro-structure of timbre and the macro-shape of musical writing by imposing temporal relations on signal descriptors. It has also been designed in order to embed symbolic and signal units in a common framework. The system relies on the previously presented time series database and audio querying paradigms that acts as an organized lexicon. This database contains symbolic information as well as signal descriptors of sounds from several collections, allowing to perform queries on their temporal shapes. This system called *spectral maquettes* for OpenMusic (OM) allows temporal and structural interactions between musical units of various nature. It is the first step towards a wider composition framework that could navigate between symbolic and spectral data.

16.4.1 Motivation

During last century, instrumental music has undergone a turnaround, tending towards to increasingly transcend traditional categories of writing (pitch, duration and intensity) and extensively embracing inharmonic, noisy and strongly time-varying sounds. The advent of technology has pushed back the frontier between noise and music, as evidenced by the approach of Pierre Schaeffer. Indeed, he considered [?] the need to "*replace the limited variety of instrumental timbres possessed by an orchestra with the infinite variety of noise timbres obtained through special mechanisms*". We can deplore the lack of compositional systems that could catch up to these new practices. Nowadays composition systems usually follow the traditional harmonic paradigm and are inscribed in a punctual time of writing. However many musical contexts (figures, textures, gestures) seem to not only work by overlaying stationary sounds but rather by carving, vertically and horizontally, the sound. From this, emerges a contemporary issue in computer music research : the *signal / symbolic* interaction. These two research streams have long been impervious to each other, partly because of the apparent heterogeneity of their

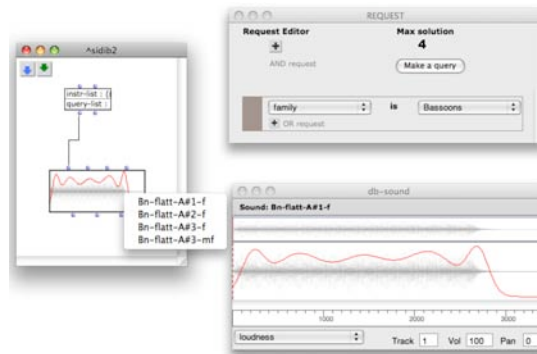


Figure 84: Make an instance of the db-sound class by a SQL query. The new db-sound instance contains four sounds. A special editor displays the current one and different descriptors of it.

objects of study. On the one hand, the analysis and synthesis of digital sounds allowed the production of sounds previously unheard. On the other hand, the algorithmic approach focused on the symbolic structures of musical notation. The composition of sound mixtures is located precisely at the crossroads of these two lines of research. If it claims to create timbres, it is through a writing process. It is therefore an essential meeting point between the symbolic and spectral domains. So it can be grasped by computer music only at the condition of convergence of the signal and the algorithmic approach.

In order to handle the complexity of sound mixtures, we propose a composition system that could effectively deal with the stylistic aspects of computer music composition focused on timbre. It has also been designed in order to embed symbolic and signal units in a common framework. We can use the knowledge obtained with the MOTS framework through the previously presented database. The innovative audio querying paradigms can then allow to perform queries based on the temporal evolution of descriptors rather than simply on static criteria over mean descriptors. This architecture thus establishes an organized lexicon for our system. In order to compose with sounds at the level of the musical discourse, we use OpenMusic's maquette system as a basis. A maquette can be seen as a "meta score" embracing in a single document musical notation and visual programs.

16.4.2 Implementation in OpenMusic

First, we defined a new class of sounds belonging to the database (db-sound). User can make instances of this class by sending SQL queries to the spectral database using the previously defined time series matching techniques (see Figure 84). More than a sound, a db-sound is a collection of sounds whose cardinality is specified by the user. As common sound files, db-sound instances are first class citizen (i.e. they can be included in visual programs, scores or other OpenMusic editors). By using db-sound instances along with the general visual programming tools of the OpenMusic environment, composers are thus provided with means to develop complete processes in relation with sound mixture issues. In this context, the functional program structures let them manage the complexity by maintaining hierarchical control from musically relevant abstractions down to synthesis processes.

In particular we have implemented a new class of maquettes called Spectral Maquettes. The maquette is an extension of the notion of visual program with additional spatial and temporal dimensions, allowing one to put the elements of the composition framework (data structures and processes) in close relation to these two dimensions. In a maquette editor, the boxes (called temporal boxes) represent functional units (programs) producing musical outputs. The position and graphical properties of these boxes are associated with a temporal and structural sense; particularly, the horizontal axis of the editor represents time, so that the position and horizontal extension can be related to offsets and durations. The temporal boxes can also be linked by functional connections so that the whole maquette may finally be considered as a program, comprising functional and temporal semantics. Temporal relations and constraints can therefore be set between the boxes by setting the temporal parameters in their corresponding patches. Hence, the calculus can determine the time structure. In addition, the possibility of embedding maquettes in a maquette allows for the construction of hierarchical temporal structures.

The spectral maquette is an extension of the classical one aiming at fomenting the relation of musical material and processes with micro-time structures. Indeed sound mixture requires a particular attention to the temporal evolution of the individual components as well as the organization, sequencing, and articulations of a global musical form. We show in Figure 3 an example of spectral maquette. When a db-sound is put into the maquette the subjacent box shows the wave form of the sound augmented with several information like clue points (e.g. attack, sustain and release) or the evolution of some descriptor (e.g. loudness, noiseness, spectral centroid, etc.). Boxes into a spectral maquette can be related between them in different ways. In Figure 3 (a) we can see how the start time of a sound is linked to the first note of a melody; moving one of these boxes force to move the other one in order to respect this constraint. In part (b) three sounds are superposed, we can see in this example several types of temporal relations: the lowest one is attached to a flag in the temporal ruler, this sound will start always at this time; the middle box, in turn, finishes the lowest one (other Allen's relation can be easily imposed between boxes); finally relations between the highest box and the middle one are defined by means of internal points of sounds (these points can be defined by hand or coming from queries on the database). In part (d) the db-sound is taken as an input of a program which builds a sequence of notes whose melodic profile follows the evolution of the loudness of the sound. Finally in part (e) a viewer allows to select a particular descriptor and calculate the result of the superposition of all db-sounds in the maquette. In a first time this viewer allows to see the global evolution of the superposition, but it is not difficult to imagine that by changing the curve the user can launch a query that search for new boxes approaching this new curve. In the previous example, we tried to show how spectral saquettes enable to structure the musical information at different levels :

- the static level of the form, i.e. the organisation of boxes in time.
- the dynamic and paradigmatic level of the form (i.e. constraints between the temporal boxes).
- the syntactical level, i.e. the calculus building the musical discourse inside the temporal boxes.
- the material level, i.e. db-sounds taking part in the maquette

These four levels of information are obviously interconnected. The most important advantage in the spectral maquette concept is to offer a visualisation of this interaction

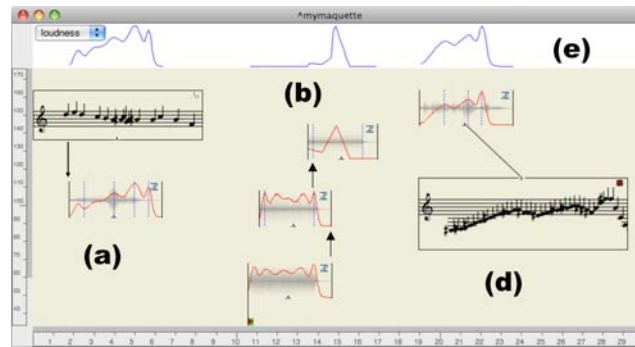


Figure 85: Concrete example of the usage of spectral maquettes. Several functional, macro and micro-temporal relations are defined between boxes.

and at the same time, interactive control of it. This produces a source of exploration and experimentation :

- Recombination at the form level. Temporal boxes are moved and stretched in time without changing the other three levels.
- Modification of functional relations. We do not change the position of the blocks but their causal relation.
- Syntax modification. Algorithms which build the material can be changed according to the compositional goals.
- Change of the material by launching new queries to the database.

When combining these procedures, sophisticated musical experiments can be performed. It is important to note that in the spectral maquettes, the user does not work directly with sounds. The resulting sound of a spectral maquette is synthesized by the server in order to create a new sound of the database or simply to be listened.

17

CONCLUSIONS OF THIS PART

Part VI

CONCLUSIONS

18

FUTURE WORK

18.1 THE MOTS PARADIGM

18.1.1 *Wider applications*

We believe that the framework of MOTS matching could lead to powerful applications. This algorithmic problem has, to the best of our knowledge, never been addressed. Hence, we are expecting to develop its application range. We already tried with the subsequent HV-MOTS *classification* to provide a maximal range of application fields. However, the idea of *matching* the time series in a multiobjective manner can also provide an interesting variety of applications. In this section, we envision some potential applications for future work.

Sound computing

As the structural choices underlying the MOTS framework were drawn from observations of our auditory perception, it would appear logical to extend its scope of audio applications. In this study, we already applied the MOTS system to the analysis of sound samples and shown that it was fit for both matching and classification. However, an interesting application would be to extend this scope to musical song matching. For example, a dual high-level analysis can lead to obtain a disjoint set of rhythmic and melodic time series for a same piece. Therefore, the multiobjective flexibility could also allow to separately match rhythmical patterns and melodic series on a larger time scale.

Medical diagnosis

The information gathered for medical analysis is usually provided by sensors that collect repeated measurements at fixed intervals of time. Therefore, they produce time series observations such as electrocardiogram (ECG), electroencephalogram (EEG), blood pressure and so forth. These signals are mostly used to monitor the state of health of patients. Hence, one of the major concern of medical informatics has been to help diagnosis and prognosis. Because of the temporal nature of medical signals (and the wealth of source where measurements could be extracted), an interesting line of work would be to study the application of the MOTS framework in these fields.

Genetic analysis

In bioinformatics, a lot of research has been devoted to the analysis and matching of DNA sequences. These sequences can be considered as time series, either in their raw format or by transforming these into random walk-type series by putting different weights on nucleic bases. Furthermore, there also exists a flourishing literature in multiobjective optimization for computational biology. Therefore, the field of genetic analysis appears as a promising topic of study for our algorithmic framework.

18.1.2 Hybrid analysis

Along this document, we studied the application of flexible notions of multidimensional similarity to sets of time series. However, an interesting line of future work could be found in the expansion of these concepts to even more complex problematics.

Combining views

First, an interesting line of work lies in experimenting some variations on the concepts of multiobjective similarity. For example, the normalization of time series data allows to separate its shape from its first statistical moments. Therefore, in the context of audio applications, it appears interesting to perform queries that jointly minimize the mean, standard deviation and temporal shape of the same feature.

Multiobjective subsequence matching

The main idea of the MOTS framework is analyze the similarity on sets of *distinct* time series. However, a very interesting study would be to apply the same concept to studying the similarity of *subsequences* of the *same* series. Therefore, this would allow to obtain several objectives that rely to the *same time series*. Hence, this study would extend the MOTS framework to work *inside* univariate time series. This relates to the stimulating question of time scales continuum, that we will detail in Section 18.8.

18.1.3 Interaction and representation

Multiobjective analysis have been applied to a wide range of problems in the past years. However, a traditional problem lies in the representation of the solutions and the possibilities of interaction offered by different systems. Indeed, multiobjective analysis can span a very high-dimensional space. Therefore, the representation of the Pareto solutions is bound to provide only a dithered view of the optimization space. Therefore, an interesting line of future work would be to study the potential solutions for representation and interaction with the MOTS framework. A good premise to this study would be to study dimensionality reduction techniques that could be applied and how these could influence the design of interaction schemes.

18.2 MOSEQ / QVI

18.2.1 Applications in audio workflow

The two audio querying paradigms introduced in Section 7.5 looks promising given the results of their user validations presented in Chapter 8. However, their use in a seldom context do not appear relevant to nowadays audio workflows. Indeed, the current music production architectures are converging towards fully integrated framework. Sometimes, these even provide video edition features or other unrelated supplements. An interesting direction of research would be to study how the MOSEQ and QVI paradigms could be integrated in a music production workflow. Furthermore, their integration in such frameworks would require a form of interaction with other musical modules. Therefore, it would be interesting to study how the inputs and outputs of these paradigms could be generalized towards natural integration.

18.2.2 Relevance feedback

As we discussed in Section 7.5, the multiobjective space embed a notion of *preference* towards different solutions. Indeed, each solution in the Pareto front is correlated with an underlying set of weights. Selecting one of these solutions rather than others implies a clear preference over different optimization dimensions. Therefore, an interesting direction of research would be to provide a weighting mechanism for each dimension in order to find relative preferences. Then, this approach could be extended to compute persistent “global” weights from one search to another. Finally, this could lead to automate this relevance mechanism. In this line of thought, an interesting user study would be to analyze if subsequent searches regularly goes towards the same optimization directions. This idea could ultimately bring an auto-modifying framework that could intelligently adapt to different users.

18.3 HV-MOTS CLASSIFICATION

18.3.1 Audio applications

The HV-MOTS classifier isn inspired by concepts from the auditory perception. Therefore, it would seem logical to apply it further in different fields of audio applications. Several topics of MIR are fundamentally directed towards classification problems and could benefit from this new classification scheme. For example, it would be interesting to apply the HV-MOTS classifier to wider time scales and elements of study. Furthermore, several classification problems are still open and would be interesting to study with the flexibility of this framework. These topics comprise *musical genre* classification, *musical mood* inference and *music composer* classification.

18.3.2 Scope of application

As we have seen in Chapter 11, because of the ubiquitous nature of time series, the HV-MOTS classification paradigm can be applied to a wide range of problems. We tried to retrieve a maximal diversity of datasets on which to analyze the accuracy of this framework. However, as put forward by Demsar [105] even the largest scale studies may not be able to project results on datasets that were not part of the original study. Therefore, an interesting direction of future work would be to analyze the statistical superiority of the HV-MOTS classifier on an even larger scale. We believe that the flexibility of this classification framework could prove itself to be statistically superior on a wide range of topics. In fact, multiobjective optimization techniques and time series analysis have a long history of successful but separate applications. Hence, an interesting line of future work would be to analyze the introduction of time series in multiobjective analysis framework. Conversely, the multiobjective flexibility could enhance traditional results in time series classification. Therefore, the HV-MOTS classifier could still be applied to a whole range of scientific fields comprising *economics*, *chemical engineering*, *process design*, *scheduling*, *bioinformatics*, *computational biology*, *climate analysis* and *medical diagnosis*

18.3.3 *Multiobjective subsequence classification*

As we argued along our study, the main property of the HV-MOTS classifier is that it is constructed to provide multidimensional similarity assessments. Therefore, it is theoretically meaningless to apply it to univariate time series classification. Indeed, in this case the hypervolume selection provides a classification decision equivalent to the 1-NN classifier. However, an interesting line of future work could arise from the study of *multiobjective subsequence classification*. A first approach to this concept would be to try to extract several objectives from a set of univariate time series related to the same feature. Therefore, it would allow to study the accuracy of the HV-MOTS paradigm in univariate time series classification. However, this requires to study how to extract those subsequences, choose the corresponding cutting points and, of course, its computational complexity.

18.4 HEART SOUNDS BIOMETRY

As we have seen in Chapter 12, the S-Features and HV-MOTS approach allow to construct an accurate identification system based on heart sounds. Even if the error rates already outstands the state-of-art proposals, a wide amount of work can still be performed to enhance these results.

18.4.1 *Segmentation procedure*

As we have seen, the segmentation procedure is one of the most critical module when analyzing the final results. It seems that the tuning this procedure may dramatically impact the error rates. Therefore, one of the most important enhancement for this system would be to provide a segmentation procedure directed towards the specificities of heart sounds. Indeed, our segmentation procedure has been selected inside a set of available methods for sound segmentation. However, more specific techniques could be developed, with a special interest towards the use of the Stockwell transform as the basis for segmentation.

18.4.2 *Features computation*

We developed a specific set of features based on the Stockwell transform (cf. Section A.2.2) for heart sounds biometry. It seems that the S-Transform is a valid choice as it allows to outperform classic tools of features analysis for this specific problem. However, this set still seems somehow limited compared to the wealth of research that has been performed in the field of audio analysis. Therefore, an interesting line of work would be to expand the set of S-Features for heart sounds biometry. Nevertheless, another interesting work would be to analyze the applicability of this feature set to other problematics. Of course, a natural choice would be to try to apply this set to audio problems. However, because of the low-frequency oriented nature of this feature set, it might also be relevant to other problems.

18.4.3 *Factors of influence*

We tried to list in Section 12.5 the factors that could influence the performances of a heart sounds biometric system. Therefore, it would be of utmost interest to study the

true impact of these factors on the performance of the system. The most interesting parameters to study would be the *physiological factors* and *heart diseases*. However, the success of such analyses requires the collection of more exhaustive datasets. For example, studying the *physiological factors* in a population of subjects imply that we should record the *same* subjects in varying physiological conditions (rest, light and hard physical exercise, sleep). Therefore, a first mandatory step would be to perform an exhaustive data collection campaign. This would allow to obtain databases that could compete with nowadays knowledge on other biometrics.

18.5 ON HEART DISEASES DETECTION

After analyzing the results of the heart sounds identification system, a natural idea may come to mind. If the system is so efficient in identifying persons through the sound their heart produce, it should easily detect heart diseases based on the same information. Indeed, as this approach is able to discriminate extremely small variations of the heart sound signals, it should be even easier to discriminate between extremely varying signals such as heart diseases. However, as for the previous section, such a study would require to collect a very wide dataset of heart diseases to support a valid statistical conclusion. Ideally, the dataset should even contain some subjects that are recorded *before* and *after* developing a heart disease. Obviously, it appears to be a very hard case to find but also implies numerous ethical issues.

18.6 ORCHESTRATION

We proposed in Section 15.4 a new system for *temporal* and *abstract* computer-aided orchestration. We have shown that it provides a powerful enhancement over existing systems. However, as we discussed in Chapter 15, musical orchestration is an extremely complex topic which encompass several open issues and problematics. We outline here some potential trends and avenues for the next years of research in musical orchestration.

18.6.1 *Signal and symbolism*

Musical orchestration offers a unique framework to study the interactions between the signal and symbolic research streams. Linking the symbolic world of writing with the realm of signal could lead to potentially powerful applications. We will discuss a hypothetic system implementing these ideas as an interesting line of future work in Section 18.7.1. We focus here on performing an analysis of such ideas through the existing body of musical works.

Relations between scales

In order to find a relation between the signal and symbolic layers, it would be interesting to analyze existing musical pieces through an automatic knowledge extraction procedure. Therefore, this analysis would focus on finding musical configurations where links between spectral patterns and symbolic writing are denoting a logical joint evolution. This study would require a coordinated analysis between the score, instruments information and resulting audio features. This joint analysis could identify the relationships between these two worlds. In order to link the acoustic time series and

symbolic objects, an interesting lead would be to use artificial intelligence technique between a set of signal features and symbolic information.

Extending to generic sound mixture

Even if classical orchestration is still an open problem, an interesting avenue of research would be to extend it to generic sound mixtures and electronics. This requires an environment design which could take into account the specificities of electronic parts. This goal is already being assessed by composers through mixed orchestration situations. In these contexts, we need to find configurations in which these two acoustic modes can blend together, without giving an impression of two dichotomic worlds. This problematic can be approached by subverting the problem of orchestration to the problem of timbre

18.6.2 *Intelligent music notation*

Current music notation systems does not support the *contextual rules* of music writing. Hence, they provide no further analysis of playing techniques and acoustic properties of individual instruments. While sophisticated word processors feature aiding tools to verify *grammatically correct* texts, no current music notation systems offer such “musical grammar” aid. Particularly in the context of orchestration, musical writing can be a painstaking task. Indeed, careful considerations of fingering, breaks in register, masking effects and so forth should be taken into consideration. Therefore, an interesting line of work would be to develop music notation systems, capable of recognizing incorrect writing and recommending alternative corrections. Therefore, a large amount of knowledge should be gathered and represented efficiently on both the symbolic and spectral aspects of instrumental capacities. This knowledge gathering procedure would require an amount of manual retrieving to ascertain the validity of information. It should also provide acoustical properties of the instrument and how these can affect questions of orchestral effects.

Automatic knowledge gathering

Based on the previous goals, a line of future work would be to automate the gathering of instrumental knowledge for orchestration. The idea would be to construct this knowledge incrementally from an automatic analysis on existing scores. Starting from a “*nothing is possible*” configuration, each apparition of a particular symbolic sequence adds it to the knowledge, and recurrence makes the sequence easier to reproduce (on a probabilistic scale with a confidence degree). Based on the symbolic scores, each sequence, chord and fingering can be added in the knowledge database and considered feasible. The more instance we found in analysis, the easier it should be to perform the related sequence. Then a cross-analysis with corresponding signal recordings and audio features could allow to obtain deeper and more interesting relationships between symbolic and signal descriptors. This approach could lead to a powerful knowledge database on signal/symbolism relations.

18.6.3 *Emergence phenomenon*

Orchestration embeds some of the most misunderstood acoustic phenomena. First, the blending of orchestral instruments can almost make it impossible to find back

these original constituents (a phenomenon termed as *orchestral fusion* in the literature). Therefore, the orchestral timbre contains complex non-linear interactions. Several phenomenon have been observed in music listening and are termed as *emergence phenomenon*, but are yet to be fully explained. First, the well-known *masking effect* in which some elements disappear from our perception, the *chorus* and *unison* effects (in which a group of the same instruments play the same notes) and the *rugosity* and *coupling* effects. It seems crucial to understand these phenomena in order to integrate them in the orchestral reasoning. Therefore, we should study some groupings of instruments which generates some unexpected elements inside the spectrum. A first solution would be to obtain recordings of several emergence phenomena and try to find the relation between their spectrum.

18.7 CLOSING THE GAP BETWEEN ALL WORLDS

As we have talked extensively in Chapter 1, orchestration lies at the crossroads of signal and symbolism but also thrives between micro-temporal evolution and macro-temporal articulations. We tried to sketch a first approach that could help in translating the intent of a composer in the process of orchestration through the *spectral maquettes*.

Hence, an interesting line of future work would be to allow a more intuitive and therefore easier control of sound mixtures. This system should provide an interaction between the spectrum (spectro-) of sounds and their evolution in time (-morphology). However, although the spectral content and temporal evolution are undoubtedly linked, we need to separate concepts to describe them. Based on these remarks, we envision a first system that could try to bridge the gap between those worlds through the concept of constraints inference.

18.7.1 Constraint inference system

The main idea behind this system would be to extract the implicit relationships that are created by a composer while writing music, as exemplified in Figure 86. For example, a chord will surely be related to its context, ie. previous notes. However, we can also find some spectral and temporal relations inside that chord. Ultimately, the goal would be to let the users compose freely and then obtain a constraint network that could be modified algorithmically. That way, we can infer some high-level relations as well as micro-level properties. By processing a simultaneous analysis of the symbolic score and the corresponding audio features evolution, we could obtain a graph of constrained relationships, explaining the link on both types of viewpoints. The key here is in combining time series inference and knowledge mining approaches. For instance, by using a spectral database, it is possible to perform a musical writing and at the same time predict the audio features of each melodic line. Therefore, the links between the harmonic rules and signal features could be explicitated through a constraint network (as shown in Figure 86) Then based on the inferred network, it would be possible to modify and then solve alternate networks. The main problem for the inference system would be to filter the knowledge induced from an automatic constraint inference. Indeed, some of the inferred constraints could be irrelevant. For example, note X is played 143,35 seconds after the beginning of the piece seems like a rather dull information. Oppositely, the temporal relations between this note and others are crucial to the final musical work.

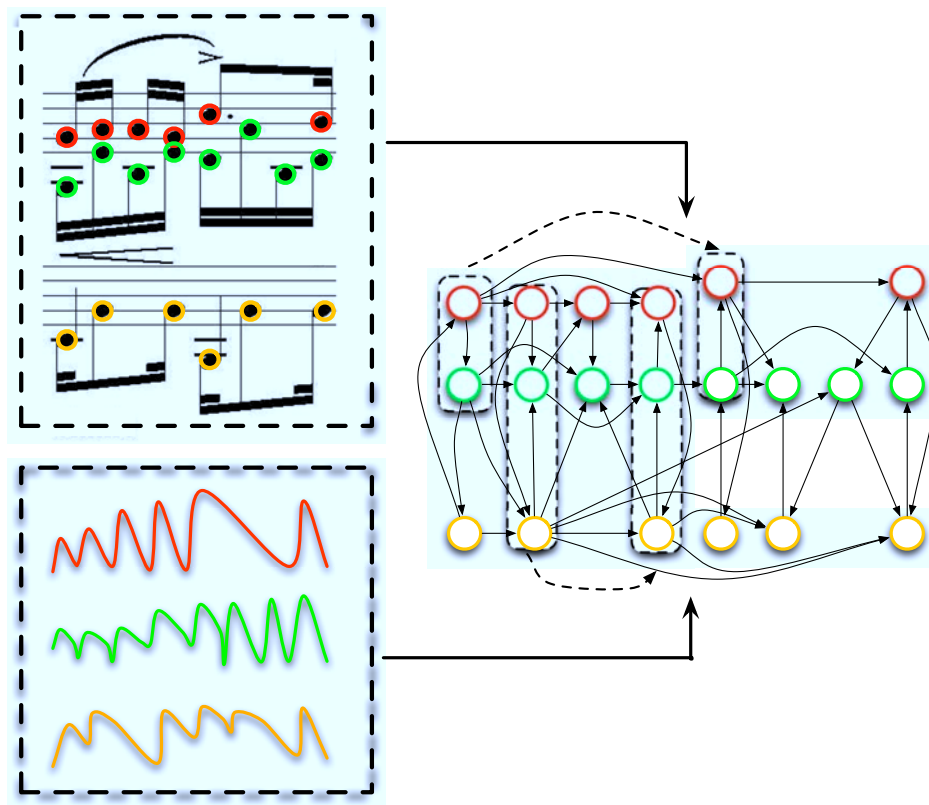


Figure 86: A complete system of constraint inference that allow to bridge the gap between the symbolic realm of musical writing and the signal world of timbre. A simultaneous analysis of the symbolic score and corresponding audio features evolution could provide a graph of constrained relationships, explaining the link on both types of viewpoints.

18.7.2 Several views on constraints

This constraint system could be articulated over several aspects of an orchestration system. Therefore, it should constitute a multi-level hierarchical system. We try to list here every beneficial module that could lead to an unifying framework of constraint handling.

- *Instrumentation and static constraints* - Symbolic constraints between instruments, notes, playing modes or other symbolic parameters are a first extensive topic of study. Instrumentation constraints on fingering, articulations and writing could exhibit links between successive orchestrations on the symbolic level. (eg. *continuity*, *diversity* and so forth). In fact, a tremendous amount of work on instrumentation is yet to be done (we detailed some research avenues in Section 18.6.2). This approach requires an extensive knowledge database with fingering, articulations speed and every instrument technical capacities.
- *Temporal constraints and macro-articulations* - We can draw several advantages from linking a purely compositional writing framework with a more generic spectral orchestration system (as ATO-MS presented in Section 15.4). This brings back the link between different temporal scales from the micro-temporal evolution of signal properties to the macro-temporal articulations of musical discourse (we will try to develop this avenue of research through the idea of *time scales continuum* in Section 18.8). In terms of constraints, this would imply to establish a “meta-level” of constraints to provide a continuous link between several temporal scales. Regarding macro-articulations a straightforward example of application would be to find *orchestral paths* between two orchestrations (examples of constraints include “*minimize the number of instruments change*” or “*maximize the variation on a set of descriptors*”)
- *Formal orchestration constraints language* - An interesting line of work would be to define an operational language for orchestration by formalizing its basic operations. This would provide a formal language for orchestration by starting with simple operations and then complexifying these by iterative refinement. Therefore, this language should be defined from its constituting elements (*musical atoms*), combined through simple operations (eg. amplification, transposition, overlay) and towards a constraint checking and solving procedure. The first step is to characterize the language primitives that could range from sinusoidal components to complete melodic lines (polymorphism of objects). The set of operations should also encompass this range between signal operations (amplification, transposition and stretching) to symbolic processing (concatenation, overlay, alignment and assembly). Finally, the operators should also embed the different time scales through a notion of *operator granularity*.

18.8 ON A TIME SCALES CONTINUUM

We have seen along this study, that we can work with widely varying time scales, each of which provide a different information focus. We try here to foresee the basis for a more generic problem. We hypothesize a more *elastic notion* of time based on a *continuum of time scales*, which could handle the notions of *temporal granularity*. This question rise from the simple fact that we can not define a clear boundary between the notions of *micro* and *macro-time*. For instance, a continuous glissando of several

seconds is considered as a micro-temporal evolution. However, in the same time frame, we can construct complete melodic lines for different instruments. These are obviously considered as a macro-temporal articulation of musical discourse. Therefore, in musical problematics at least, there seems to be a discrepancy of the classical notion of time scales. Based on the idea of scales continuum, it would be interesting to extend our reasoning to a dense set of temporal domains. This would allow to perform a similar reasoning at every time scale. Hence, processing each temporal granularity in the same manner would provide a wider homogeneity. We discuss this problematic through the prism of orchestration but it could easily be transposed to any type of temporal process.

18.8.1 *From micro to macro*

The connection between micro and macro temporalities would allow to incorporate timbral structure in the compositional act. Hence, we must find ways to reassemble the temporal reasoning from the smallest levels all the way up to macro-temporal level. A first step would be to study the states of interaction between different temporalities and how the micro-time can influence the macro-level. This leads to a "*problem of uncertainty*" in which the relationships between local and global representations of time-varying spectrum are yet to be connected. A solution could be to use of a bottom-up pattern recognition process based on small collections of primitives, that could result in clustering of features into patterns. It is also interesting to consider the creative part of parsing in its connective ability (and not merely reducing) that can encompass several layers of description. This approach would seek to organize these layers into an integrated musical structure with multiple facets. It is unlikely that this task is looking for a 'one-to-one' correspondence between the layers. It should rather recognize that the musical context emerges from nonlinear interactions. Dimensional reduction can give a good angle to attack this problem. It would allow to analyze the evolution of musical structure and how this movement is correlated with its spectral rendering. There may be no linear continuity between the different time scales. However, the solution can not consist of the outright closure of their borders, but should consider all the possibilities of interaction and hence articulation between these different scales. This question of time granularity could connect the treatment of events from different time scales in the same encompassing framework.

18.8.2 *Macro-temporal articulations*

The second step of this problematic lies in defining the evolution of sound mixtures over large period of time. As we discussed earlier, the temporal constraints can define paths between several orchestral states. An interesting possibility would be to provide symbolic paths and arbitrarily choose the temporal evolution of a spectral dimension to be optimized. Therefore this study should try to find the link that can unite the spectral flow of two orchestrations in a temporal manner. Therefore, it would be interesting to study long-term modes for "articulated" targets whose timbre is changing continuously.

We have shown along this study that a musical problematic can give birth to powerful analysis schemes, far beyond their own study. By questioning the nature of artistic reasoning and at the same time drawing inspiration from our musical perception and gaining insights from these mechanisms to drive our choices of algorithms, we have been able to create novel and powerful approaches for generic querying and classification.

First, we introduced the problem of *MultiObjective Time Series* (MOTS) matching and its formalization. We discussed the core differences between this novel framework and its multivariate counterpart, and have shown how it could lead to more flexible sets of retrieval solutions. We introduced two efficient algorithms to solve this problem by relying on the concept of *approximate hyperplane*. We have analyzed their relative merits on synthetic and real datasets, which exhibited that the *hyperplane* algorithm was able to solve the MOTS problem in sub-linear time with respect to both database and objective cardinalities. However, further analysis of the results show that the variance of querying time are more important with this method as its complexity depends on the distribution of data. Based on the MOTS framework, we were able to formalize innovative audio querying by introducing two new querying problematics in the field of audio samples retrieval. The *MultiObjective Spectral Evolution Query* (MOSEQ) in which users can directly draw the temporal evolution of audio features and *Query by Vocal Imitation* (QVI) which allow users to perform a vocal imitation of the sound they are seeking. Both paradigms are a direct application of the flexibility introduced by the MOTS matching problem.

We then performed extensive and thorough user studies on the temporal evolution of audio features. These studies allowed to validate the hypotheses on which the MOTS framework, which exhibited the concept of *directions of listening*. We have shown that complex audio features can be perceived by this perception is influenced by the temporal evolution of other features. We have further shown the consistency of these directions of listening in generic similarity tasks. Finally, an extensive usability evaluation of the MOSEQ and QVI querying paradigms have shown that they are indeed interesting and intuitive querying paradigms for audio retrieval.

Based on these results, we extend our scope of study and show how to apply these notions of flexible similarity evaluation to classification problems by introducing the HyperVolume-MOTS (HV-MOTS) classification scheme. We show that even within the multiobjective framework that avoids merging distances into a single measure, we can still rank classes by relying on the hypervolume dominated by each. We discuss the relation between this novel classification framework and other distance-based classifiers but also to more generic classification schemes and further provide a discussion on its main advantages and drawbacks. We provide a large scale study of the performances of this classification technique on a wide range of datasets that covers several scientific fields. We show the statistical superiority of the HV-MOTS classifier over well-established classification schemes and over state-of-art results on the same datasets. Based on the HV-MOTS classifier, we show how to construct a biometric identification system for heart beat sounds by once again considering listening as an art and developing a specific set of features based on the Stockwell transform, called

S-Features. We show that using heart sounds as a biometric feature provide a reliable identification and that this feature is not affected by the phenomenon of *template ageing* over a time span of two years, supported by the recordings collected in the Mars 500 isolation study. We show the application of the HV-MOTS framework to audio problems through generic audio samples classification and sound morphology.

We show how this knowledge gained through broader applications can be put to use in the field of musical orchestration. We introduce a new orchestration system based on a algorithm that use the MOTS framework and that rely on a entropic segmentation procedure. We show that this new algorithm outperforms the previous approaches for computer-aided orchestration. We then present other artistic applications of the MOTS framework.

On a more epistemological level, we can see that the analysis of musical problematics offers a very powerful framework to study wider scientific topics. As we have discussed along this study and in our directions of future work, the inherent temporal nature of music writing exhibit very stimulating problematics. Through the connections that exist between the micro-level of spectral properties and the macro-time of musical discourse, rise the question of *temporal granularity* and even further the idea of a *dense temporal scales continuum*.

Part VII

APPENDIX

A

APPENDIX

A.1 HV-MOTS DATASETS DESCRIPTION

ARABIC-DIGIT (UCI)

Field	Spoken digit recognition
Summary	This dataset contains the recordings of 10 spoken arabic digit by 88 speakers. Each digit is repeated 10 times for each speaker. The data comes from 44 males and 44 females native Arabic speakers.
Features	Sound files have been analyzed to obtain 13 Mel-Frequency Cepstrum Coefficients (MFCCs) time series.
Classes	Spoken arabic digits ([0 – 9])
Samples	8800 samples (10 digits x 10 repetitions x 88 speakers)
Sampling	Computation from sound files at 11025 Hz sampling rate, 16 bits with Hamming window. A pre-emphasis filter was applied to the original signal.
Source	This dataset is part of the UCI repository [131] and is extensively described in [165, 164].
Results	A Vector Quantization (VQ) with a Maximum Weight Spanning Tree (MWST) leads to a final mean classification accuracy of 93.12% over every classes with single classes results varying from 85.55% to 99.00%

ARTIFICIAL CHARACTERS (UCI)

Field	Character recognition
Summary	This dataset has been artificially generated by using first order theory to describe the structure of ten capital letters of English alphabet. Each instance is described by a set of segments which imitate the way an automatic program would segment an image.
Features	Each segment is represented by X and Y values for the starting and ending points.
Classes	10-classes problem representing the capital letters A, C, D, E, F, G, H, L, P and R
Samples	6000 samples (600 for each class)
Source	This dataset is part of the UCI repository [131] and is described in Botta et al. [53]
Results	A Genetic Algorithm coupled with an histogram local optimization leads to a recognition rate of 98.68%

AUSTRALIAN-SIGNS (UCI)

Field	Sign recognition
Summary	This dataset consists of samples of Auslan (Australian Sign Language) signs that were recorded with multiple sensors on a powered glove. Examples of 95 signs were collected from five signers with a total of 6650 sign samples
Features	The glove recorded 10 different time series feature for each sign as the x , y and z position of the hand, the <i>roll</i> , <i>pitch</i> and <i>yaw</i> of the hand orientation and the bending for <i>thumb</i> , <i>forefinger</i> , <i>index</i> , <i>ring</i> and <i>little</i> fingers.
Classes	Each sign represents a different class which amounts to 95 different classes.
Samples	6650 samples (varies for each class)
Source	This dataset is part of the UCI dataset repository [131] and is extensively described in [199].
Results	For the low quality set (Nintendo Power-glove), best results are obtained by a Hidden Markov Model (HMM) classifier with 71.2% classification accuracy [200]

AUSTRALIAN-SIGNS-HQ (UCI)

Field	Sign recognition
Summary	This dataset is a highest quality version of the Australian-signs dataset. However, recordings were made with a single signer. It contains 27 examples of each of the same 95 signs for a total of 2565 signs collected from a native signer using high quality position trackers on both hands.
Features	The same feature set is used, however this time both hands were recorded simultaneously.
Samples	2565 samples (27 for each class)
Articles	This dataset is part of the UCI dataset repository [131] and is described in [200].
Results	For this dataset, best results are obtained by the Naive Segmentation based on TClass algorithm with a 94.5% accuracy [200].

BCIII-01-TUBINGEN

Field	Brain-Computer Interface
Summary	This dataset assess the motor imagery in ECoG recordings, for a set of imaginary movements. Cued motor imagery (left pinky, tongue) from one subject; training and test data are ECoG recordings from two different sessions with about one week in between
Features	The recordings contains 64 ECoG channels (0.016-300Hz).
Classes	2-class discrimination between motor imagery of <i>left pinky</i> and <i>tongue</i>
Samples	378 samples (278 <i>training</i> and 100 <i>test</i>)

Sampling	1000Hz sampling rate
Source	This system is part of the 2003 BCI Competition III Blankertz et al. [49] and is described in Lal et al. [232]
Results	The use of three specific features extracted with a Common Spatial Subspace Decomposition (CSSD) exhibited a 91% accuracy Quiguo et al. 313.

BCIIII-02-ALBANY

Field	Brain-Computer Interface
Summary	The goal of this dataset is to estimate to which letter of a 6-by-6 matrix with successively intensified rows and columns the subject was paying attention to (P300 speller paradigm)
Features	Recordings contains 64 EEG channels (0.1-60Hz)
Classes	36 classes
Samples	185 samples (85 training and 100 test)
Sampling	240Hz sampling rate
Source	This dataset is part of the 2003 BCI Competition III Blankertz et al. [49] and is described in Farwell and Donchin [125].
Results	A mixture of Gaussian Kernel SVMs allow to obtain a classification accuracy of 73.5% for 5 trials Rakotomamonjy and Guigue [318].

BCIIII-03A-GRAZ

Field	Brain-Computer Interface
Summary	This dataset is composed of EEG recordings of cued motor imagery with 4 classes (<i>left hand, right hand, foot, tongue</i>) from 3 subjects (ranging from quite good to fair performance). Performance is measured using the kappa-coefficient.
Features	60 EEG channels (1-50Hz)
Classes	4-class problem between <i>left hand, right hand, foot</i> and <i>tongue</i>
Samples	840 samples (70 trials per subject for each class)
Sampling	250Hz sampling rate
Source	This dataset is part of the 2003 BCI Competition 2003 Blankertz et al. [49].
Results	A multi-class CSP based on Fisher ratios obtained a kappa coefficient of 0.7926.

BCIIII-03B-GRAZ

Field	Brain-Computer Interface
Summary	Motor imagery with non-stationarity problem and online feedback (non-stationary classifier) with 2 classes (left hand, right hand) from 3 subjects.
Features	2 bipolar EEG channels 0.5-30Hz
Classes	2 class problem between <i>left</i> and <i>right hand</i>
Samples	2760 samples (460 trials per subject for each class)
Technical	125Hz sampling rate
Source	This dataset is part of the 2003 BCI Competition 2003 Blankertz et al. [49]. The performance measure is maximal steepness ($t_0=3s$) of mutual information [bits/s],
Results	Results provided in Lemm et al. [237] shows that a probabilistic selection of causal Morlet wavelets bins allows a maximum accuracy of 89.3 %

BCIIII-04A-BERLIN

Field	Brain-Computer Interface
Summary	This dataset is composed of cued motor imagery with 2 classes (<i>right hand</i> , <i>foot</i>) from 5 subjects. For 2 subjects most trials are labelled (resp. 80% and 60%), while from the other 3 less training data are given (resp. 30%, 20% and 10%). The challenge is to make a good classification even from little training data, thereby maybe using information from other subjects with many labelled trials.
Features	118 EEG channels (0.05-200Hz)
Classes	2 classes (<i>right hand</i> and <i>foot</i>).
Samples	1400 samples (280 trials per subject)
Technical	1000Hz sampling rate
Source	This dataset is part of the 2003 BCI Competition III Blankertz et al. [49] and described extensively in Dornhege et al. [113].
Results	The results provided in Wang et al. [396] shows that a Neural Network trained on features extracted from a Common SubSpace Decomposition (CSSD) obtains a classification accuracy of 92.98%.

BCIIV-01-BERLIN

Field	Brain-Computer Interface
Summary	Motor imagery for an uncued EEG classifier application (for <i>hand</i> and <i>foot</i>); evaluation data is a continuous EEG which also contains periods of idle state
Features	64 EEG channels (0.05-200Hz)

Classes	3 classes (<i>hand</i> , <i>foot</i> and <i>idle</i> state)
Samples	1400 samples (200 trials per subject)
Sampling	1000Hz sampling rate
Source	This dataset is part of the BCI Competition IV Blankertz et al. [50]. The performance measure is the mean squared error with respect to the target vector
Results	The output of the competition shows that a clustering procedure applied on a Principal Components Analysis (PCA) allows to obtain a MSE of 0.382.

BCIIV-03-FREIBURG

Field	Brain-Computer Interface
Summary	Hand movement directions in MEG. The data set contains directionally modulated low-frequency MEG activity that was recorded while subjects performed wrist movements in four different directions.
Features	10 MEG channels (filtered to 0.5-100Hz) located above the motor areas
Classes	4 classes (<i>forward</i> , <i>backward</i> , <i>left</i> and <i>right</i> wrist movements)
Samples	480 samples (120 for each classes with 60 trials per subject)
Technical	The trials were cut to contain data from 0.4 s before to 0.6 s after movement onset and the signals were band pass filtered (0.5 to 100 Hz) and resampled at 400 Hz.
Source	This dataset is part of the BCI Competition IV Blankertz et al. [50].
Results	The final obtained accuracy is 46.9% over the 2 subjects with 59.5% for the first and 34.3% for the second. Feature extraction and reduction is then fed to a Genetic Algorithm (GA) which decide the features to use for classification with a combination of linear Support Vector Machine (SVM) and Linear Discriminant Analysis (LDA)

BIOMAG-2010

Field	EEG analysis
Summary	The goal of this dataset is to detect whether subjects are attending to the left or right visual field on each trial based on the MagnetoEncephaloGram (MEG) of the subjects.
Features	274 MEG channels recorded independently
Classes	2-class problem (between <i>left</i> and <i>right</i> visual field).
Samples	780 samples
Source	This dataset is described in Van Gerven and Jensen [382]
Results	The results reported in Van Gerven and Jensen [382] shows that a Support Vector Machine (SVM) algorithm provide an accuracy of up to 69% correctly classified trials.

BRAIN-COMPUTER

Domain	Brain-Computer Interface
Summary	This dataset contains ElectroEncephaloGram (EEG) recordings that were acquired while the subject performed 80 left hand movement imaginations and 80 right hand movement imaginations
Features	The recordings covers 2 channels of EEG (C ₃ and C ₄), it also contains the ElectroOculoGram (EOG) and surface ElectroMyoGram (EMG) from the extensor muscles of both hands.
Classes	2-class problem between <i>left</i> and <i>right</i> hand movement imagination
Samples	160 samples (80 per class)
Technical	128 Hz sampling rate with a 1 / 50μV calibration.
Source	This dataset is described extensively in Coyle et al. [101], Guger et al. [157].
Results	Classification accuracy of 89.28% (S ₁), 88.13% (S ₂) and 93.13% (S ₃) is reported, which leads to a mean classification accuracy of 90.18%. However these results use different sets of features for each subject. The true best classification accuracy over the same feature set is 85.75 % using an Adaptive AutoRegressive (AAR) model with a Linear Discriminant Analysis (LDA).

CHALLENGE-2011

Field	Cardiology
Summary	This dataset is composed of 12-lead ECG recordings, with the aim of providing useful feedback on the quality of the ECG signals by classifying them depending on their quality (from poor to excellent).
Features	This dataset is composed of 12 leads (<i>I</i> , <i>II</i> , <i>III</i> , <i>aVR</i> , <i>aVL</i> , <i>aVF</i> , <i>V1</i> , <i>V2</i> and <i>V3</i>)
Classes	2-class task between <i>acceptable</i> and <i>unacceptable</i> ECG recording depending on their qualities.
Samples	2000 samples (1000 for each class)
Technical	The leads are recorded simultaneously for a minimum of 10 seconds; each lead is sampled at 500 Hz with a 16-bit resolution.
Source	This dataset is part of the <i>PhysioBank</i> archive [146] posted as the 2011 challenge.
Results	The challenge was intended to see the selectivity and specificity of algorithms. The best system reports a 85.9% accuracy Xia et al. [407] using the spectrum radius of a matrix of regularity.

CHARACTER-TRAJECTORIES (UCI)

Field	Character recognition
Summary	This dataset is composed of labelled samples of pen tip trajectories recorded for individual characters writing. All samples are from the same writer, for the purposes of primitive extraction. Only characters with a single pen-down segment were considered
Features	The dataset is made of three time series for each instance, which represents the x and y position of the pen and the <i>pen tip force</i> .
Classes	20-class problem over different characters
Samples	2858 samples (varies for each class)
Technical	Recordings have been made at 200 Hz sampling rate. Data has been numerically differentiated and Gaussian smoothed, with a sigma value of 2.
Source	This dataset is part of the UCI repository [131] and is described in Williams et al. [402, 403].
Results	Reported classification accuracy of 87% in Zafar et al. [424] and 93.67% in Perina et al. [297] using a Gaussian Mixture Model (GMM) with a Hidden Markov Model (HMM).

DACHSTEIN

Domain	High altitude medicine
Summary	Each instance represents EEG and ECG data for one cardiac cycle that were acquired at 900 m and at 2700m altitude. The subject performed a reaction time task. The data shows the influence of the loss of oxygen on event-related desynchronization (ERD) and event-related synchronization (ERS) and heart rate variability.
Features	The recordings covers 2 channels of EEG (C3 and C4) and the ECG recordings.
Classes	2-class problem between 900m and 2700m recordings
Samples	698 samples (324 at 900m and 374 at 2700m)
Technical	256 Hz sampling rate with a 1 μ V calibration
Source	This dataset is described in Guger et al. [158] where analysis is performed to show the influence of loss of oxygen
Results	- No classification accuracy reported -

DIGIT-HANDS (UCI)

Field	Digit recognition
Summary	This dataset contains 250 samples from 44 writers performing each of the ten digits in random order.
Features	The tablet sends x and y coordinates and <i>pressure</i> level values of the pen at fixed time intervals.
Classes	10-class problem with each digit ($[0 - 9]$)
Samples	10992 samples (varies for each class)
Technical	The sampling rate is 100 miliseconds
Source	This dataset is part of the UCI repository [131] and is described in Alimoglu and Alpaydin [9].
Results	Classification accuracy of 97.8% is achieved with a 3 – NN classifier.

EEG-ALCOHOLISM (UCI)

Field	Medical analysis
Summary	This data comes from a large study to examine EEG correlates of genetic predisposition to alcoholism. It contains measurements from 64 electrodes placed on subject's scalps.
Features	Each recording is composed of 64 EEG electrodes time series

EEG-alcoholism-2

Classes	2-class problem between <i>alcoholic</i> and <i>control</i> subjects
---------	--

EEG-alcoholism-3

Classes	3-class problem between S_1 <i>object</i> , S_2 <i>match</i> and S_2 <i>non-match</i> situations
---------	--

EEG-alcoholism-6

Classes	6-class problem for S_1 <i>object</i> , S_2 <i>match</i> and S_2 <i>non-match</i> situations between <i>alcoholic</i> and <i>control</i> subjects
Samples	650 samples
Technical	Data was recorded at 256 Hz sampling rate
Source	This dataset is part of the UCI repository [131] and is described in [429].
Results	The results provided in Zhong and Ghosh [431] show that a Multivariate HMM allows a classification accuracy of 90.5%. However, this result is obtained with only 10 pre-selected measurements. The true classification accuracy is 78.5%

EEG-EPFL

Field Human-Machine Interaction

Summary This dataset is composed of EEG recordings from 9 subjects. Experimental protocol shown 6 images on a computer screen and subjects were ordered to focus on a single image during a session. Images were then randomly flashed the same number of times each. The goal is to find which image the subject was focusing on.

Features The raw EEG recordings are collected from 34 electrodes. The ordering of electrodes is *Fp1, AF3, F7, F3, FC1, FC5, T7, C3, CP1, CP5, P7, P3, Pz, PO3, O1, Oz, O2, PO4, P4, P8, CP6, CP2, C4, T8, FC6, FC2, F4, F8, AF4, Fp2, Fz, Cz, MA1* and *MA2*. This amounts to 34 time series features for each trial.

Eeg-epfl-7

Classes 7-class problem between all subjects depending on either *target* (1-6) and *non-target*.

Eeg-epfl-12

Classes 12-class problem between all subjects depending on either *target* (1-6) and *non-target* (1-6)

Eeg-epfl-36

Classes 36-class problem formed by all combinations between the *target* (1-6) and *flashing image* (1-6)

Samples 26646 samples

Technical Recordings were performed at a 2048 Hz sampling rate.

Source The dataset and experimental protocol is described in Hoffmann et al. [179]

Results The classification accuracy is separate for each subject and attains between 95 and 100% for every subject using a Bayesian Linear Discriminant Analysis (BLDA).

FORTE

Domain Climatology (lightning prediction)

Summary This dataset is split into three different problems depending on the number of classes they contain. Each dataset is aimed at predicting the type of lightning observed through recordings of the power density.

Features Ground density of Electro-Magnetic Power (EMP) recording

Samples 121 instances in each dataset.

Source The dataset is extensively described in [103].

Results Classification results of 77.5% accuracy are reported in Bernecker et al. [42] using a Shared-Nearest-Neighbors (SNN) algorithm with the Longest Common SubSequence (LCSS) distance.

Forte-2class

Classes 2-class problem where distinction should be made between *cloud-to-ground* and *intra-cloud* lightning.

Forte-6class

Classes 6-class problem where distinction should be made between CG (Positive-Initial Return Stroke), SR (Subsequent Negative Return Stroke), IR (Negative Initial Return Stroke), I (Impulsive Event), I₂ (Impulsive Event Pair) and KM (Gradual Intra-Cloud Stroke) lightning events

Forte-7class

Classes 7-class problem where distinction should be made between CG (Positive-Initial Return Stroke), SR (Subsequent Negative Return Stroke), IR (Negative Initial Return Stroke), I (Impulsive Event), I₂ (Impulsive Event Pair), KM (Gradual Intra-Cloud Stroke) and O (Off-Record) lightning events

GAITPDB

Field Medical analysis

Summary This database contains measures of gait from patients with idiopathic Parkinson Disease (PD) (mean age: 66.3 years; 63% men), and healthy control subject (mean age: 66.3 years; 55% men). A disturbed gait is a common, debilitating symptom of PD; patients with severe gait disturbances are prone to falls and may lose their functional independence. The database includes the vertical ground reaction force records of subjects as they walked at their usual, self-selected pace for approximately 2 minutes on level ground.

Features This dataset was recorded using 8 sensors under each foot, which gives the *vertical ground reaction force* for each foot as well as the *total force* under each foot. This amounts to 18 time series features.

Classes 2-class problem between *healthy* and *parkinson* subjects based on gait disturbances.

Samples 306 samples (214 parkinson and 92 healthy subjects)

Technical Recordings were made at 100 Hz sampling rate

Source This dataset is part of the PhysioBank database [146] and is described in Frenkel-Toledo et al. [133]

Results Results introduced in Lee and Lim [236] shows that Neural Network with weighted fuzzy membership functions on Wavelet Transform coefficients allow to obtain a 77.33% classification accuracy.

HANDWRITTEN (UCI)

Field	Character recognition
Summary	This dataset contains 8235 online handwritten assamese characters. The “online” process involves capturing the data in real-time while the characters are written on a digitizing tablet with an electronic pen. This dataset was collected from 45 writers, each of which contributed 183 recordings
Features	The acquisition program records the handwriting as a stream of X and Y coordinate points using the appropriate pen position sensor along with the pen-up and pen-down switching. No pressure level was recorded.
Classes	This is a 183-classes problem. Each writer contributed 52 basic characters, 10 numerals and 121 assamese conjunct consonants.
Samples	8235 samples (45 for each class).
Source	This dataset is part of the UCI repository [131]
Results	A classification accuracy of 92% with HMM to 96% with SVM is reported but only for digits (so only a 10-classes problem)

IONOSPHERE (UCI)

Field	Radar analysis
Summary	This dataset was collected by a radar system in Goose Bay, Labrador. The targets were free electrons in the ionosphere. "Good" radar returns are those showing evidence of some type of structure in the ionosphere. "Bad" returns are those that do not; their signals pass through the ionosphere.
Features	There were 17 pulse numbers for the Goose Bay system. Instances in this database are described by 2 attributes per pulse number, corresponding to the complex values returned by the function resulting from the complex electromagnetic signal which amounts to 34 features.
Classes	2-class binary task between <i>good</i> and <i>bad</i> radar returns.
Samples	358 samples
Technical	The recording system consists of a phased array of 16 high-frequency antennas with a total transmitted power of 6.4 kilowatts. Received signals are processed using an autocorrelation function whose arguments are the time of a pulse and its number.
Source	This dataset is part of the UCI repository [131] and is described in Sigillito et al. [352]
Results	The best results is a 94.2 % classification accuracy reported in Eggermont et al. [118] using a Genetic Programming (GP) approach.

JAPANESE-VOWELS (UCI)

Field	Speaker identification
Summary	This dataset covers the topic of speaker identification for nine male speakers which uttered two Japanese vowels successively.
Features	Each instance is represented by 12 LPC cepstrum coefficients time series.
Classes	9-class problem representing each speaker
Samples	640 samples (varies for each class)
Technical	10 kHz sampling rate analyzed with a 25.6ms window size and a 6.4ms hop size.
Source	This dataset is part of the UCI repository [131] and is described in Kudo et al. [229].
Results	The proposed classifier exhibits a classification accuracy of 94.1%, while a 5-state continuous Hidden Markov Model attained up to 96.2%

LIBRAS (UCI)

Field	Movement recognition
Summary	The dataset contains 15 classes of 24 instances each. Each class represents to a hand movement type in LIBRAS (official brazilian signal language).
Features	In each frame, the centroid pixels of the segmented objects (the hands) are found. These successive points compose the discrete version of the curve which is represented by the x and y position time series. All curves are normalized in the unitary space.
Classes	15-class problem for <i>swings</i> , <i>arcs</i> , <i>circle</i> , <i>lines</i> , <i>zigzag</i> , <i>waves</i> , <i>curves</i> and <i>tremble</i> movements
Samples	360 samples (24 for each class).
Technical	In the video pre-processing, a time normalization is carried out selecting 45 frames from each video, by following an uniform distribution.
Source	This dataset is part of the UCI repository [131], described originally in Dias et al. [108] and later used in Schliebs et al. [341]
Results	An evolving Spiking Neural Network (eSNN) allows to obtain a 88.59% classification accuracy.

MEG-MIND-READING

Field	MEG analysis (mind reading)
Summary	This dataset has been gathered by recording MagnetoEncephaloGraphy (MEG) signals from subjects while looking at video clips belonging to three different categories (<i>artificial</i> , <i>football</i> and <i>nature</i>) and two long recordings from a <i>Bean</i> and <i>Chaplin</i> film. The goal is to perform a mind reading algorithm which can predict the category of the video clip based on the MEG recordings.

Features	Signals are recorded with a 306-channel MEG system. The final features are 204 <i>planar gradiometer channels</i> .
Classes	5-class dataset divided into <i>artificial, football, nature, bean</i> and <i>chaplin</i> categories.
Samples	1330 samples (677 training and 653 testing)
Technical	Signals are divided in one-second periods, digitized at a 1000Hz sampling rate. Signals have been low-pass filtered to 50Hz and downsampled to a sampling rate of 200Hz.
Source	This dataset has been assembled for the ICANN'2011 conference Klami et al. [221] and studied in Huttunen et al. [184]
Results	The proposed Bayesian Cross-Correlation Analysis (CCA) classifier provide a 68.0% classification accuracy.

NEUROTUCHO-SOCIAL

Field	Monkey ECoG analysis
Summary	ECoG data and eye position were recorded. The monkey is confronted to different social situations with or without human presence. The three situations were neutral, frighten and threaten.
Features	The ECoG recordings have been performed over 128 different channels
Classes	8-class interactions (<i>frighten-alone, frighten-human, neutral-alone, neutral-human, threaten-alone, threaten-human</i>)
Samples	300 samples
Technical	Monkey was sitting with head fixed. His arm motion was also restrained. ECoG data were sampled at 1KHz and unit is micro volt. Motion data and eye tracking data were sampled at 120Hz. Start and stop point of all data were synchronized.
Source	This dataset is part of the Neurotycho Nagasaka et al. [279] collection.
Results	- No classification results reported -

NEUROTUCHO-VISUAL

Field	Monkey ECoG analysis
Summary	ECoG data and eye position were recorded. There was a monitor in front of the monkey. Grating pattern that moves in eight direction was presented on the screen. There was no fixation required.
Features	ECoG recordings with 128 channels
Classes	8-type visual stimulus classification task that represents each <i>directions</i> of the <i>grating patterns</i> .
Samples	200 samples

Technical	Monkey was sitting with head fixed. His arm motion was also restrained. One cycle of sinusoid pattern was 27mm with speed at 108mm/sec (4Hz). Distance between monkey and screen was 490mm. Blank and stimulus pattern were switched alternatively every 2 sec. ECoG data were sampled at 1KHz and unit is micro volt. Motion data and eye tracking data were sampled at 120Hz. Start and stop point of all data were synchronized.
Source	This dataset is part of the Neurotycho Nagasaka et al. [279] collection.
Results	- No classification results reported -

PEN-CHARACTERS (UCI)

Field	Character recognition
Summary	This dataset is composed of upper and lowercase characters, digits, and other spanish characters, for a total of 62 different characters collected from 11 different writers which performed 2 repetitions.
Features	Each character is represented by a sequence of segments summarized by their <i>X</i> and <i>Y coordinates</i> for each point.
Classes	62-classes problem.
Samples	1364 samples (2 repetitions for 11 writers for each class)
Source	This dataset is part of the UCI repository [131] and is described in Prat et al. [307]
Results	Classification accuracy is shown to be up to 89.15% if using a DTW matching algorithm on segment-based representation with a NN-rule.

PEN-CHARACTERS-2 (UCI)

Field	Character recognition
Summary	This dataset is composed of ASCII and non-ASCII characters which amount to a total of 97 different characters collected from 60 different writers which performed 2 repetitions.
Features	Each character is represented by a sequence of segments summarized by their <i>X</i> and <i>Y coordinates</i> for each point.
Classes	97-classes problem.
Samples	11640 samples (2 repetitions for 60 writers for each class)
Source	This dataset is part of the UCI repository [131] and is described in Castro-Bleda et al. [80]
Results	Classification accuracy is presented in Castro-Bleda et al. [80] of 91.8% classification accuracy if using a template matching algorithm with a NN-rule however this results is <i>only for a restricted set of 62 characters</i> .

PERSON-ACTIVITY (UCI)

Field	Movement analysis
Summary	People used for recording of the data were wearing four tags (ankle left, ankle right, belt and chest). Each instance is a localization data for one of the tags. The goal of this dataset was to detect falls only but we also test the accuracy in classifying all movement classes.
Features	The four tags worn by subjects sent the x , y and z position for <i>left</i> and <i>right</i> ankle, <i>chest</i> and <i>belt</i> which amounts to 12 distinct features.
Classes	This dataset features instances of different postures and actions with <i>walking</i> , <i>falling</i> , <i>lying</i> , <i>lying down</i> , <i>sitting</i> , <i>sitting down</i> , <i>standing up from lying</i> , <i>on all fours</i> , <i>sitting on the ground</i> , <i>standing up from sitting</i> and <i>standing up from sitting on the ground</i> .
Samples	164860 samples (varies for each class)
Source	This dataset is part of the UCI repository [131] and is described in [203]
Results	Classification accuracy is shown to be 72% for machine learning agents, 88% for expert-knowledge agents and 91.3% for meta-prediction agents. However, this results are for detecting falls only.

PHYSICAL-ACTION (UCI)

Field	Movement analysis
Summary	Three male and one female subjects (age 25 to 30), who have experienced aggression in scenarios such as physical fighting, took part in the experiment. Throughout 20 individual experiments, each subject had to perform ten normal and ten aggressive activities.
Features	The data acquisition process involved eight skin-surface electrodes placed on the upper arms (biceps and triceps), and upper legs (thighs and hamstrings), which corresponds to eight input time series for all muscle channels (ch1-8). Each time series contains around 10000 samples (15 actions per experimental session for each subject). The electrodes were placed on <i>right bicep</i> (C1), <i>right tricep</i> (C2), <i>left bicep</i> (C3), <i>left tricep</i> (C4), <i>right thigh</i> (C5), <i>right hamstring</i> (C6), <i>left thigh</i> (C7) and <i>left hamstring</i> (C8)

Physical-action-2

Classes	2-class problem between <i>aggressive</i> and <i>normal</i> actions
---------	---

Physical-action-20

Classes	20-class task divided between <i>elbowing</i> , <i>frontkicking</i> , <i>hammering</i> , <i>headering</i> , <i>kneeing</i> , <i>pulling</i> , <i>punching</i> , <i>pushing</i> , <i>sidekicking</i> , <i>slapping</i> , <i>bowing</i> , <i>clapping</i> , <i>handshaking</i> , <i>hugging</i> , <i>jumping</i> , <i>running</i> , <i>seating</i> , <i>standing</i> , <i>walking</i> and <i>waving</i> actions.
Samples	80 samples (4 for each class)

Technical	The subjects' performance has been recorded by the Delsys EMG apparatus, interfacing human activity with myoelectrical contractions.
Source	This dataset is part of the UCI repository [131]
Results	A Genetic Programming (GP) algorithm allows to obtain 73.3% classification accuracy Theodoridis and Hu [376]. However these results are for six of the actions only.

PTBDB

Field	Cardiology
Summary	The ECGs were collected from healthy volunteers and patients with different heart diseases
Features	Each record includes 15 simultaneously measured cardiac signals. The conventional 12 ECG-leads (<i>i, ii, iii, avr, avl, avf, v1, v2, v3, v4, v5, v6</i>) together with the 3 Frank lead ECGs (<i>vx, vy, vz</i>).
Classes	9 separate class representing healthy patients and patients suffering different heart conditions. The classes are therefore <i>healthy controls, myocardial infarction, cardiomyopathy, bundle branch block, dysrhythmia, myocardial hypertrophy, valvular heart disease, myocarditis</i> and <i>miscellaneous</i> .
Samples	2750 instances from 549 records performed on 290 subjects (aged 17 to 87, mean 57.2; 209 men, mean age 55.5, and 81 women, mean age 61.6; ages were not recorded for 1 female and 14 male subjects). Each subject is represented by one to five records.
Technical	Each signal is digitized at 1000 samples per second, with 16 bit resolution over a range of ± 16.384 mV. with $0.5 \mu\text{V}/\text{LSB}$ (2000 A/D units per mV).
Source	This dataset is part of the PhysioBank database [146] and is described in [54].
Results	The results provided in Gudmundsson et al. [156] shows that a Random Forest (RF) classifier provide a 75.1% classification accuracy

Ptdb-1

Technical	This dataset contains single heart beats for classification
-----------	---

Ptdb-2

Technical	This dataset contains two heart beats in each instance
-----------	--

Ptdb-5

Technical	This dataset contains five heart beats in each instance
-----------	---

ROBOT-FAILURES (UCI)

Field	Robotics
Summary	This dataset contains force and torque measurements on a robot after failure detection. Each failure is characterized by 15 force/torque samples collected at regular time intervals
Features	The robots measure a set of x , y and z <i>force position</i> as well as a x , y and z <i>torque</i> measured after failure, which amounts to a total of 6 individual time series features.
Classes	This dataset is divided into five sub-sets, each of them defining a different learning problem
Technical	Each failure instance is characterized in terms of 15 force/torque samples collected at regular time intervals starting immediately after failure detection. A total observation window of 315ms is used for each failure instance.
Source	This dataset is part of the UCI repository [131] and described [67].
Results	A set of five feature transformation strategies (based on statistical summary features, discrete Fourier transform, etc.) allows to obtain a maximal classification accuracy of 80%.

Robot-failures-lp1

Classes	Failures in approach to grasp position 24% normal - 19% collision - 18% front collision - 39% obstruction
Samples	88

Robot-failures-lp2

Classes	Failures in transfer of a part 43% normal - 13% front collision - 15% back collision - 11% collision to the right - 19% collision to the left
Samples	47

Robot-failures-lp3

Classes	Position of part after a transfer failure 43% normal - 19% slightly moved - 32% moved - 6% lost
Samples	47

Robot-failures-lp4

Classes	Failures in approach to ungrasp position 21% normal - 62% collision - 18% obstruction
Samples	117

Robot-failures-lp5

Classes	Failures in motion with part 27% normal - 16% bottom collision - 13% bottom obstruction - 29% collision in part - 16% collision in tool
Samples	164

SLPDB

Field	Sleep apnea analysis
-------	----------------------

Summary The MIT-BIH Polysomnographic Database is a collection of recordings of multiple physiologic signals during sleep. Subjects were monitored in Boston's Beth Israel Hospital Sleep Laboratory for evaluation of chronic obstructive sleep apnea syndrome, and to test the effects of constant positive airway pressure (CPAP), a standard therapeutic intervention that usually prevents or substantially reduces airway obstruction in these subjects.

Features All recordings include an *ECG signal*, an *invasive blood pressure* signal (measured using a catheter in the radial artery), an *EEG signal*, and a *respiration signal* (in most cases, from a nasal thermistor). Seven-channel recordings also include a *respiratory effort signal* derived by inductance plethysmography, an *EOG signal* and an *EMG signal* (from the chin). Therefore the dataset is divided depending on the number of available features.

Classes Several tags can be applied at the same time for a subject. The first kind depends on the different sleep stage, depending on if the subject is *awake*, in *sleep stage 1, 2, 3, 4* or *REM sleep*. Then phases of apnea are recorded between different types, with *hypopnea*, *obstructive apnea* and *central apnea* that can be *with* or *without arousal*. Final different movements are recorded for *legs* and *arousal*.

Slpdb4-Apnea

Summary 2-class task between *apnea* and *normal* based on 4 features (*ECG, BP, EEG* and *Respiration*)

Slpdb4-Full

Summary 17-class task between all available tags based on 4 features (*ECG, BP, EEG* and *Respiration*)

Slpdb4-Sleep

Summary 7-class task between different sleep stage (*Wake, 1, 2, 3, 4, REM* and *Movement*) based on 4 features (*ECG, BP, EEG* and *Respiration*)

Slpdb7-Apnea

Summary 2-class task between *apnea* and *normal* based on 7 features (*ECG, BP, EEG, Respiration, Respiratory effort, EOG* and *EMG*)

Slpdb7-Full

Summary 17-class task between all available tags based on 7 features (*ECG, BP, EEG, Respiration, Respiratory effort, EOG and EMG*)

Slpdb7-Sleep

Summary 7-class task between different sleep stage (*Wake, 1, 2, 3, 4, REM and Movement*) based on 7 features (*ECG, BP, EEG, Respiration, Respiratory effort, EOG and EMG*)

Samples 4085 samples

Source This dataset is part of the PhysioBank database [146] and described in [185].

Results Between 83.24% to 88.97% classification accuracy Bsoul et al. [62] using a Multi-Scale Support Vector Classifier (MS-SVM)

SONAR

Domain Sonar analysis

Summary This dataset contains the patterns of sonar signals bouncing off either a metal cylinder or rocks. In both case, the angles and conditions varies. The goal is to correctly identify the metal cylinder returns.

Features Set of energy in 60 frequency bands over a certain period of time

Classes 2-class problem between *mine* and *rock* sonar signals.

Samples 208 samples (111 mines and 97 rocks)

Source TThis dataset is part of the UCI repository [131] and is described in **Tan and Dowe** [371].

Results The results reports a 76% classification accuracy by using a Minimum Message Length (MML) Oblique Tree.

SYNEMP

Domain Climatology (lightning prediction)

Summary The classification tasks are aimed at studying varying speed of leading edge for different classes of lightning

Features Synthetic density of Electro-Magnetic Power (EMP) recording

Classes 2-class problem between *slow* and *fast* leading edge

Samples 20000 samples (10000 for each class)

Source This dataset is described in [196].

Results - No classification results reported -

VFDB

Field	Cardiology
Summary	This database includes 22 half-hour ECG recordings of subjects who experienced episodes of sustained ventricular tachycardia, ventricular flutter, and ventricular fibrillation.
Features	Two ECG recordings based on separate electrodes.
Classes	15-class problem between different ventricular fibrillations comprising <i>atrial fibrillation, asystole, ventricular bigeminy, first degree heart block, high grade ventricular ectopic activity, normal sinus rhythm, nodal ("AV junctional") rhythm, noise, pacemaker (paced rhythm), sinus bradycardia, supraventricular tachyarrhythmia, ventricular escape rhythm, ventricular fibrillation, ventricular flutter and ventricular tachycardia.</i>
Samples	600 samples (40 per class)
Source	This dataset is part of the PhysioBank database [146] and is described in [153]
Results	A band-pass digital filtration and ECG peak detection algorithm allows to obtain a 91.5% classification algorithm in Krasteva and Jekova [228].

VICON-PHYSICAL (UCI)

Field	Physiological analysis
Summary	This dataset includes 10 normal and 10 aggressive physical actions based on various human activities. The data have been collected by 10 subjects using the Vicon 3D tracker.
Features	The <i>x, y</i> and <i>z</i> positions define the 3D position of each marker in space for <i>left</i> and <i>right wrist, elbow, ankle</i> and <i>knee</i> which amounts to a total of 26 time series features for each action.
Samples	2000 instances (10 for each action)
Technical	The duration of each action was approximately 10 seconds per subject, which corresponds to a time series of 3000 samples, with sampling frequency of 300Hz.
Source	This dataset is part of the UCI repository [131] and is described in [376].
Results	95.4% classification accuracy is reported using a Dynamic Neural Network, however it is only applied on 9 actions out of the 20 classes.

Physical-action-2

Classes	2-class problem between <i>aggressive</i> and <i>normal</i> actions
---------	---

Physical-action-20

Classes	20-class task to make distinction between <i>elbowing, frontkicking, hamering, headering, kneeling, pulling, punching, pushing, sidekicking, slapping, bowing, clapping, handshaking, hugging, jumping, running, seating, standing, walking</i> and <i>waving</i> actions.
---------	--

WALL-ROBOT

Domain	Robotics
Summary	This dataset contains ultrasound readings for a wall-following task for robotic navigations.
Features	24 ultrasounds readings and 4 minimum sensor readings.
Classes	4-class problem between <i>forward</i> , <i>slight-right</i> , <i>sharp-right</i> and <i>slight-left</i> movements
Samples	5460 samples (1365 for each class)
Technical	Sensor readings are sampled at a rate of 9 samples per second.
Source	This dataset is part of the UCI repository [131] and is described in Freire et al. [132].
Results	The results exhibit a classification accuracy of 95.58% if using a Polynomial kernel SVM.

A.2 UNICITY OF HEART SOUNDS

A.2.1 Cardiac auscultation

Auscultation

Stethoscopic auscultation still plays today a leading role in medical diagnosis. Developed with the work of Laennec in 1816 [231], the auscultation with stethoscopes has remained almost unchanged since. The late 1940s has witnessed the arrival of a new graphical method for cardiac analysis: the *PhonoCardioGram* (PCG) which associated with the recording of pulsatile phenomena allowed a precise diagnosis of cardiac pathologies. Despite its successful development, this technique has been gradually replaced since the 1970s by ultrasound imaging obtained by continuous and pulsed Doppler. However, the required echocardiographic equipment is still extremely expensive. Furthermore, the extensive training required for recording and interpretation of ultrasounds makes it a highly specialized technique. Auscultation, meanwhile, has the merit of simplicity and can be performed anywhere at any time. The last years have also witnessed the emergence of electronic stethoscopes directly equipped with noise reduction and enhancement filters. Auscultation has kept his sense today as a unique diagnostic tool for the detection of cardiac anomalies. PCG has considerably expanded its scope with the advent of graphic reproduction and signal processing techniques. For decades, it has still not be supplanted as a basis for accurate diagnosis of heart disease and is now systematically associated with the pulsatile carotid recording, or even direct recording of impulse on the chest wall.

Auscultation is usually performed through different beaches, defined by their position on the chest or abdomen. Auscultation beaches can allow doctors to hear several slightly varying cardiac cycles which enhance diagnosis by a better localisation of the anomaly. Cardiac auscultation beaches are named according to their relative positions to the heart valves. Some diseases of the thoracic aorta can lead to use a back auscultation, but most data is collected on the anterior chest wall at the aortic, pulmonary, mitral and tricuspid beaches. It is recognized that 4G2 (fourth left intercostal space

2 cm from the midline) is a beach which often provides a good summary of cardiac cycles.

Pathologies

Several cardiovascular diseases can alter the mechanisms of the human heart and consequently the sounds that it produces. We list here what kind of modifications can occur in PCG recordings and briefly discuss the related pathologies.

SUPERNUMARY NOISES In addition to S₁ and S₂, a third (S₃) and fourth (S₄) heart sounds may also be audible. These sounds occur in the diastole phase and are the result of cardiac insufficiencies. S₃ usually succeeds S₂ (from 80 to 140ms afterwards) and can be traced to a massive spontaneous filling of the left ventricle. This sound is always the sign of a ventricular pathology. S₄ precedes S₁ and is related to the occurrence of atrial contraction caused by an active ventricular filling. This sound sometimes appear in healthy adolescents but is distinctive of a pathology after 35 years of age.

SYSTOLIC NOISES Sometimes, an isolated high-frequency systolic noise occurs shortly after S₁. This noise is caused by a deficient opening of the aortic valve and is also known as *systolic ejection click*. Another opening snap can appear shortly after S₂ with the prolapse of the mitral valve. These kind of sounds are always distinctive features of heart deficiencies.

DIASTOLIC NOISES These sounds similar to a snap, originates in the mitral opening with rigidified valves. Therefore it appears approximately 80ms after S₂. Supernumerary diastolic noises are usually associated with valvular or pericardial disease. It is sometimes followed by a roll which exhibit a left ventricular filling through a restricted orifice with high atrial pressure. A high-frequency click can also occur at the beginning of the diastole because of a constricted pericardium.

SOUNDS SPLITS The usual S₁ and S₂ sounds can sometimes be heard as duplicated, which leads to “split sounds”. These are the most common anomalies and their occurrence indicates a deficiency in the right heart that causes a loss of the usual synchronism between the left and right ventricles.

HEART MURMURS Heart murmurs are usually long noises that can be heard at different times of the cardiac cycle. Murmurs are a consequence of a change in the streaming flow inside the heart. Some of these can be louder depending on their origin. Ejection murmurs are related to transvalvular flow and are very common situations where they are termed as “innocents” as they are not pathological indicators. Their classification is therefore frequently based on their timing of occurrence which can be diastolic, systolic or even continuous.

PCG Biometry

The idea of using heart sounds (*PhonoCardioGram* (PCG)) as a biometric feature has first been introduced by Phua et al. [301]. They proposed to study heart sounds through a Short-Term Discrete Fourier Transform (STDFT) in order to obtain their spectral decomposition. The decomposition is then processed with a bandpass filter in order to remove frequency bins outside the useful cardiac sounds information. Finally,

dimension reduction and spike removing allows to filter out noise and artifacts from the spectral information. In order to match the identification templates, they compared the use of Vector-Quantization (VQ) and Gaussian Mixture Model (GMM) with different sets of features (Mel-Frequency Cepstral Coefficients (MFCC) and Linear Frequency Band Cepstra (LFBC)). They extended their work in [302] with a larger database of 128 people and showed that the LFBC features computed on 256ms frames with a 4-component GMM classification scheme was the most accurate system. However, it should be noted that these works use the *same recordings* that are split in enrollment and identification sequences, therefore it cannot account for the inherent variability between different recordings conditions. Beritelli and Serrano [36] independently proposed a biometric identification system by using frequency analysis of heart sound recordings. In this paper, the authors first segmented the cardiac recordings into the two main sound components. They then applied the Chirp z-Transform (CZT) on segmented sounds to obtain the spectral information from each cardiac cycle. In order to assess the identity, feature vectors are compared to stored templates by using the Euclidean distance. They further developed this system in [37] where they used Linear Frequency Cepstrum Coefficients (LFCC) and a feature specifically designed for heart sounds called the First-to-Second Ratio (FSR) which defines the proportionality of average powers between the first (S_1) and second (S_2) heart sounds.

For the past years these ideas have simply been extended by different approaches. Tran et al. [378] investigated the use of various feature sets by exploring temporal, spectral, harmonic and rhythmic features. They further applied a feature selection algorithm in order to classify heart sounds with automatically selected features. El-Bendary et al. [121] extended the spectral decomposition approach by using the Discrete Wavelet Transform (DWT) on which they extracted correlation and cepstral features. In order to classify heart sounds they compared Mean Square Error (MSE) and k-Nearest Neighbors (kNN) classifiers and showed that the latter improved classification accuracy on their database. Jasper and Othman [195] also used the DWT with Daubechies wavelet, limited to sub-bands between 30 and 140Hz. After pre-processing and filtering the signal, they showed the superiority of the Shannon energy envelopgram on their dataset.

Other approaches have also been undertaken. Fatemian et al. [126] tried to combine ECG and PCG-based approaches for identification in order to improve classification accuracy. They processed the heart sounds with a STFT and Mel-filterbank before applying a Linear Discriminant Analysis (LDA) which allows to obtain low-dimensional vectors on which Euclidean distance can be applied. Beritelli and Spadaccini [38] proposed a statistical approach for PCG-based biometry where they use a Gaussian Mixture Model with Universal Background Model (GMM-UBM) recognition technique. This technique appears to improve accuracy over previous structural approaches.

A.2.2 *S-Features*

In order to precisely separate the properties of heart sounds, we need to use a spectral decomposition that can fit and enhance their unique characteristics. Several signal processing tools have been developed over the past years, such as the Short-Time Fourier Transform (STFT). However heart sounds have an extremely narrow useful bandwidth. Furthermore, most of their energy lie in very low frequency ranges, below the resolution power of usual decompositions. The Stockwell transform was originally developed for analyzing geophysical data [361]. This time-frequency distribution is

defined as a generalization of both the Short-time Fourier transform and the Continuous Wavelet Transform (CWT) but overcomes some of their limitations. There are several ways to express the S transform. As stated by [361], the S-Transform can be defined as a CWT with a gaussian mother wavelet multiplied by a phase factor. Therefore, the S-Transform of a function $h(t)$ is defined as

$$S_x(t, f) = \int_{-\infty}^{\infty} h(\tau) \frac{|f|}{\sqrt{2\pi}} e^{-\frac{(t-\tau)^2 f^2}{2}} e^{-i2\pi f\tau} d\tau \quad (\text{A.1})$$

Compared to classical spectral decomposition, the S-Transform provides a frequency dependent resolution which leads to a finer definition. In our case, what is more relevant is that the S-Transform (as the CWT) can give access to an extremely fine resolution even at very low frequencies (below 50Hz) which is the most useful bandwidth of heart sound signals. Therefore it allows a better distinction of spectral components. Furthermore, unlike the CWT, modulation sinusoids are fixed with respect to the time axis. This localizes dilations and translations and thus provides the same temporal resolution for each frequency bin. A fast S-Transform algorithm was proposed [61] which strongly reduces its computational complexity and therefore makes it usable in real-life applications.

Based on the computation of the S-Transform $S_x(t, f)$, we can access precise information on the frequency distributions. However, we still need to extract high-level information that could put forward the uniqueness of each heart. To that end, we developed a specifically-tailored set of features inspired by the work in musical analysis and content-based audio retrieval. We call this set of high-level information the *S-Features*.

A.2.3 Statistical moments

We can study the evolution of the shape of a distribution by computing its statistical moments. That way, we can gain insights on the distribution of the energy in various frequency bins over time.

S-CENTROID The first moment is the geometric center (*centroid*) of the distribution, it allows to measure which frequency defines the central position of the energy distribution

$$S_{\text{centroid}}(t) = \int f \cdot S_x(t, f) df \quad (\text{A.2})$$

S-SPREAD The S-Spread is based on the second statistical moment (or *variance*) which exhibits the dispersion of energy distribution over the frequency bins

$$S_{\text{spread}}(t) = \int (f - S_{\text{centroid}}(t))^2 S_x(t, f) df \quad (\text{A.3})$$

S-SKEWNESS The S-Skewness is based on the third statistical moment which measures the symmetry of the distribution. More precisely, the skewness exhibit the lack of symmetry for a distribution, where a positive value indicates that the distribution

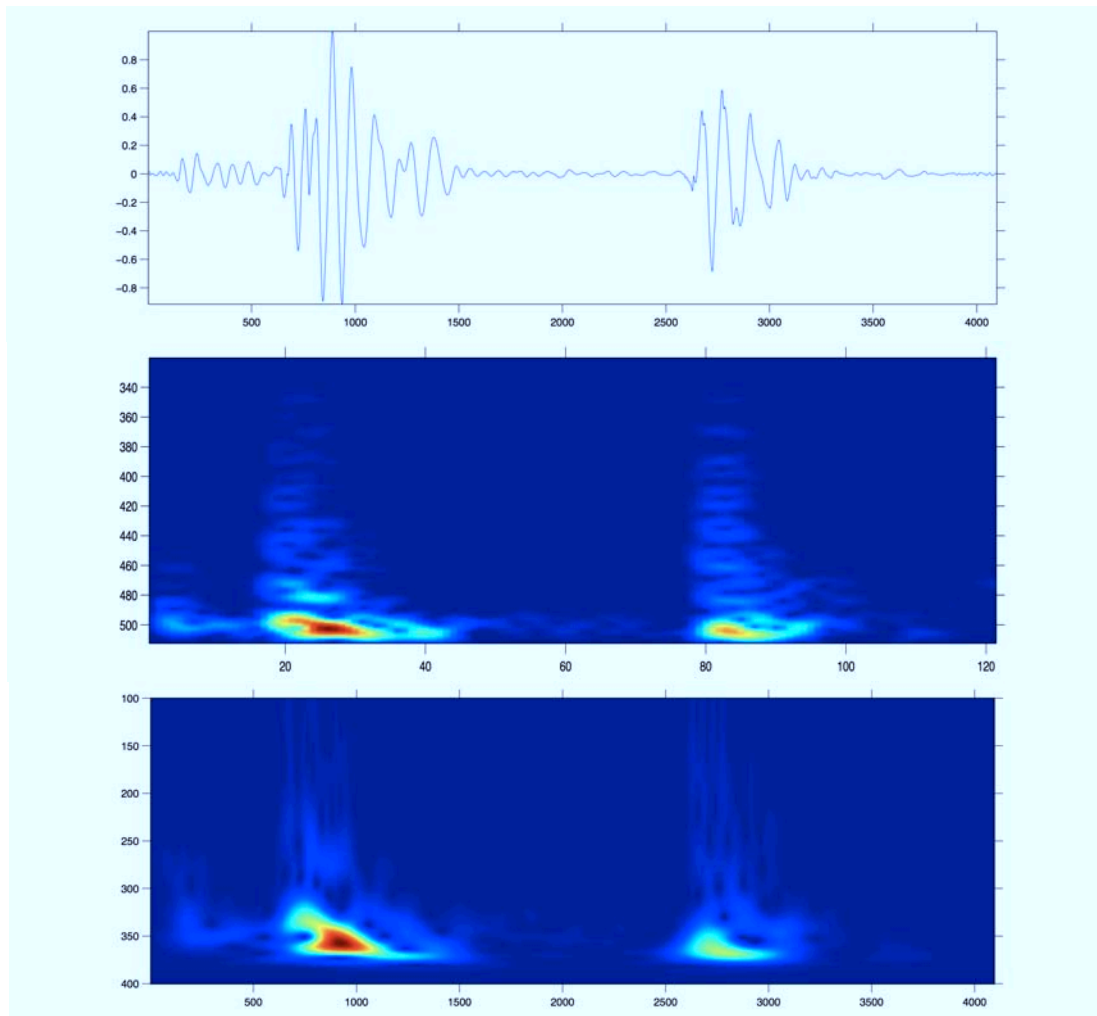


Figure 87: Comparing the resolution power of the FFT to the S-Transform for a single cardiac cycle.

have more values larger than the mean. Convertly, a negatively skewed distribution has more values lower than the mean. A symmetrical distribution has a skewness of zero.

$$S_{\text{Skewness}}(t) = \frac{\int (f - S_{\text{centroid}}(t))^3 S_x(t, f) df}{(\sqrt{S_{\text{spread}}(t)})^3} \quad (\text{A.4})$$

S-KURTOSIS The S-Kurtosis is based on the fourth statistical moment and is computed by using its corresponding definition

$$S_{\text{kurtosis}}(t) = \frac{\int (f - S_{\text{centroid}}(t))^4 S_x(t, f) df}{S_{\text{spread}}(t)^2} \quad (\text{A.5})$$

The kurtosis summarizes whether the distribution is more peaked or flat than the normal distribution.

A.2.4 Energy distribution

S-BRIGHTNESS The S-Brightness allows to study the distribution of energy over the frequency range. The idea is to fix a cut-off frequency f_c and measure the percentage of energy above that threshold, therefore yielding the temporal function

$$S_{\text{Brightness}}(t) = \frac{\sum_{f_i > f_c} S_x(t, f_i)}{\sum_i S_x(t, f_i)} \quad (\text{A.6})$$

In our implementation we tested cutoff frequencies in the set $\{65, 100, 135, 170, 200, 250\}$ Hz. We also combined these functions to provide a multi-dimensional per band descriptor called S-BrightnessBands.

S-ENTROPY The Shannon entropy, extensively used in information theory, gives a generic description of shapes. It especially allows to determine whether a distribution contains predominant peaks or not. In order to compute the S-Entropy, we use the definition of relative entropy which is independent of the sequence length

$$S_{\text{entropy}}(t) = - \frac{\sum_{i=1}^N S_x(t, f_i) \log_b S_x(t, f_i)}{\log_b(n)} \quad (\text{A.7})$$

S-FLATNESS The S-Flatness is a measure which indicates whether the distribution of energy over frequencies is smooth or if spikes are present in a spectral frame. It is computed by dividing the geometric mean by the arithmetic mean of each frame of the S-Transform

$$S_{\text{flatness}}(t) = \frac{\sqrt[N]{\prod_{i=1}^N S_x(t, f_i)}}{\left(\frac{\sum_{i=1}^N S_x(t, f_i)}{N} \right)} \quad (\text{A.8})$$

s-RMS The Root-Mean-Square (RMS) energy of each frame of the S-Transform can be computed by taking the root average of the square of each frequency bin

$$S_{\text{RMS}}(t) = \sqrt{\frac{1}{N} \sum_{i=1}^N (S_x(t, f_i))^2} \quad (\text{A.9})$$

s-ROLLOFF The S-Rolloff allows to study the temporal evolution of energy concentration. The idea is to find the frequency f_c such that a certain fraction of the total energy thresh_e is contained below that frequency.

$$S_{\text{rolloff}}^{\text{thresh}_e}(t) = f_c \mid \sum_{i=1}^{f_c} S_x(t, f_i) < t_e \quad (\text{A.10})$$

It is interesting to note that $S_{\text{rolloff}}^{0.5}(t)$ gives an approximation of the S-Centroid. In our implementation we compute the S-Rolloff for $\{65, 75, 85, 95\}$ % of energy and also provide the multidimensional S-RolloffBands.

s-FLUX The S-Flux describes the difference between two successive frames of the S-Transform. It allows to obtain a measure of “novelty” in the S-Transform by showing whether two successive frames are similar or not

$$S_{\text{flux}}(t_i) = \sqrt{\sum_{j=1}^N (S_x(t_{i+1}, f_j) - S_x(t_i, f_j))^2} \quad (\text{A.11})$$

A.2.5 Peaks distribution

The following descriptors are all based on the spectral peaks found in the S-Transform. Therefore, we first apply a peak tracking algorithm to the spectral frames. We then consider that we have access to two ordered list \mathcal{P}_i^f and \mathcal{P}_i^a which contains respectively the frequencies and amplitudes of the prohiminent peaks found in the spectrum.

s-IRREGULARITY The S-Irregularity allows to exhibit the degree of variation of contiguous peaks found in the spectrum. It is therefore obtained by computing the mean difference of amplitude between each successive pair of peaks.

$$S_{\text{Irregularity}}(t) = \frac{\sum_{i=1}^{N_p-1} (\mathcal{P}_i^a - \mathcal{P}_{i+1}^a)^2}{\sum_{i=1}^{N_p} (\mathcal{P}_i^a)^2} \quad (\text{A.12})$$

s-ROUGNESS The S-Roughness allows to study the closeness of prohiminent peaks in each frame of the S-Transform. Originally based on the auditory concept of beating phenomenon, we simplified its computation by using the average distance in frequency between all pair of peaks.

$$S_{\text{Roughness}}(t) = \frac{\sum_{i=1}^{N_p-1} \sum_{j=i}^{N_p-1} (\mathcal{P}_i^f - \mathcal{P}_j^f)^2}{\sum_{i=1}^{N_p} (\mathcal{P}_i^f)^2} \quad (\text{A.13})$$

A.3 EXTENDED HEARTS BIOMETRY ANALYSIS

We provide in this section an extended analysis of the hearts biometric system. As several parts of the proposed system may have an influence on the overall performance, we study each of them separately. These comparisons are performed following the same order as our algorithmic workflow (cf. Figure 64).

A.3.1 Pre-processing

The first step in the treatment of heart beats is to prepare the signal in order to remove its eventual impurities. As every audio signals, heart recordings are subject to *noise* and *spikes* that we pre-process thanks to different filtering techniques.

Despiking filter

The despiking filter allows to remove the *spikes* and *jitters* that might appear due to defects in the recording system. Figure 88 exhibits the influence of the despiking filter through the ROC curves of *with* and *without* use of the filter. As we can see, the use of the despiking algorithm strongly enhance the performances of the algorithm. However, this improvement seems to be more noticeable on the ROC curve (and therefore on the overall accuracy of the system) than on the rank identification rates. Therefore, the despiking filter enhance the overall performances but less on a per-person basis. This can be explained by the fact that beats which exhibit strong outliers (spikes) are “fixed” by the filter. Without filtering, these beats provoke wide error anomalies (their scores are extremely lower than normal beats). Hence, without these, the *global* score analysis exhibit the improvement. Therefore, this seems to indicate that the despiking filter is an efficient pre-processing step. However, this also indicates that the set of beats which requires such analysis seems to be of small cardinality.

Wavelet denoising

The wavelet denoising allows to handle the presence of background noise in the signal by removing its eventual presence and thus improve the *Signal-to-Noise Ratio* (SNR). When using such denoising algorithm, the choice of the wavelet is crucial to its success. Figure 89 exhibits the influence of the wavelet choice through the ROC curves of using the *Coiflets* or *Daubechies* wavelets. As we can see on the ROC curves, the use of the denoising filter also enhance the overall accuracy of the system. It seems that the use of either wavelets provide the same improvement (the two curves are almost merged). However, compared to the despiking filter, the wavelet processing seems to also strongly impact the per-person scores (as seen in the rank identification rates). It seems logical as this time, it is not a set of *single beats* that exhibit defects (throughout every person), but rather a set of *complete recordings* only for a few person. These recordings being more noisy than others, can provoke the complete scores of one individual to exhibit lower accuracy. Therefore, when these recordings are fixed, the corresponding person obtains a better ranking.

A.3.2 Segmentation

The segmentation procedure is applied after all pre-processing has been performed, in order to obtain a set of coherent heart beats from the entire recording. It is therefore

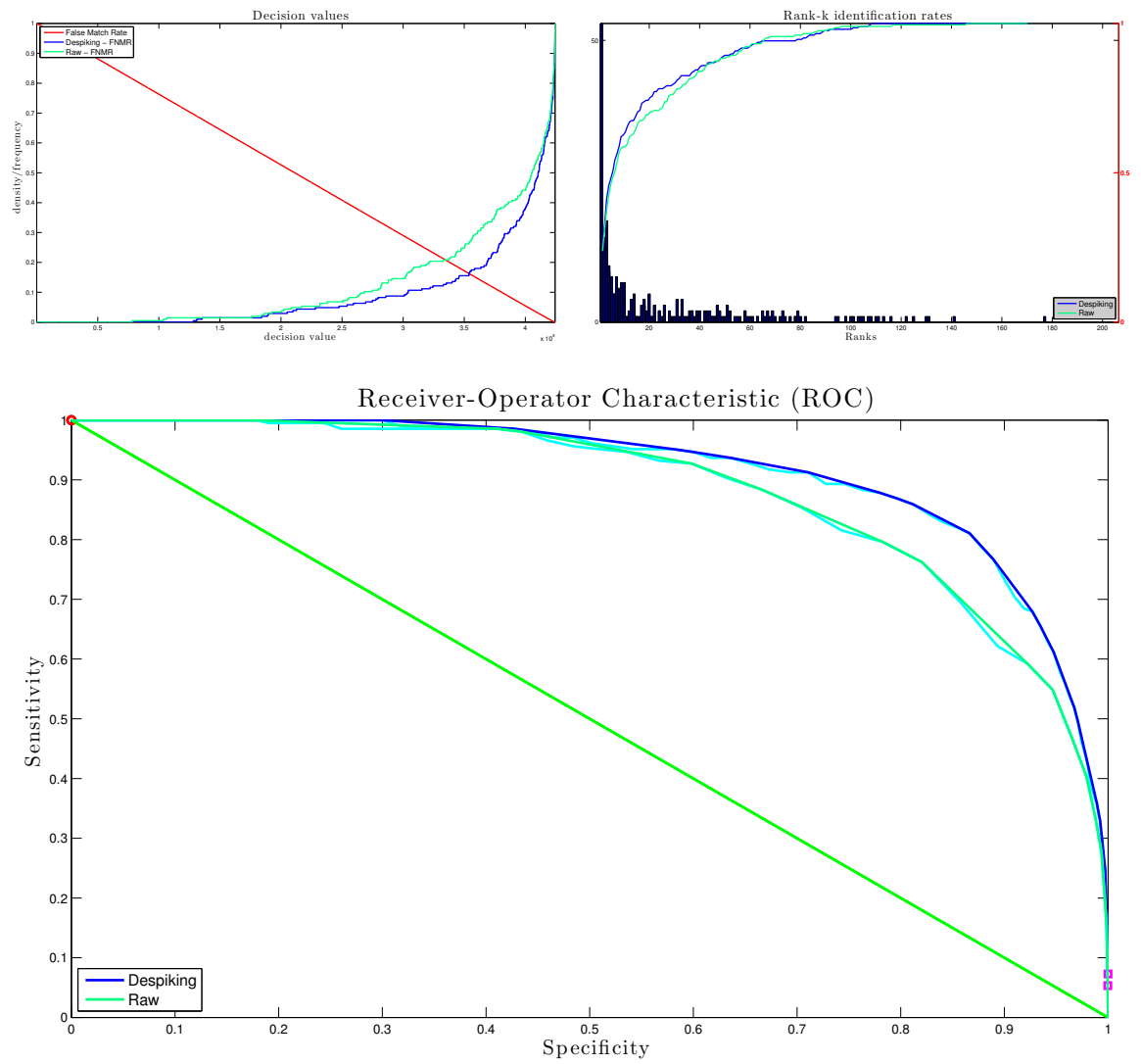


Figure 88: Influence of the despiking filter exhibited through the ROC curves of *with* and *without* use of the filter.

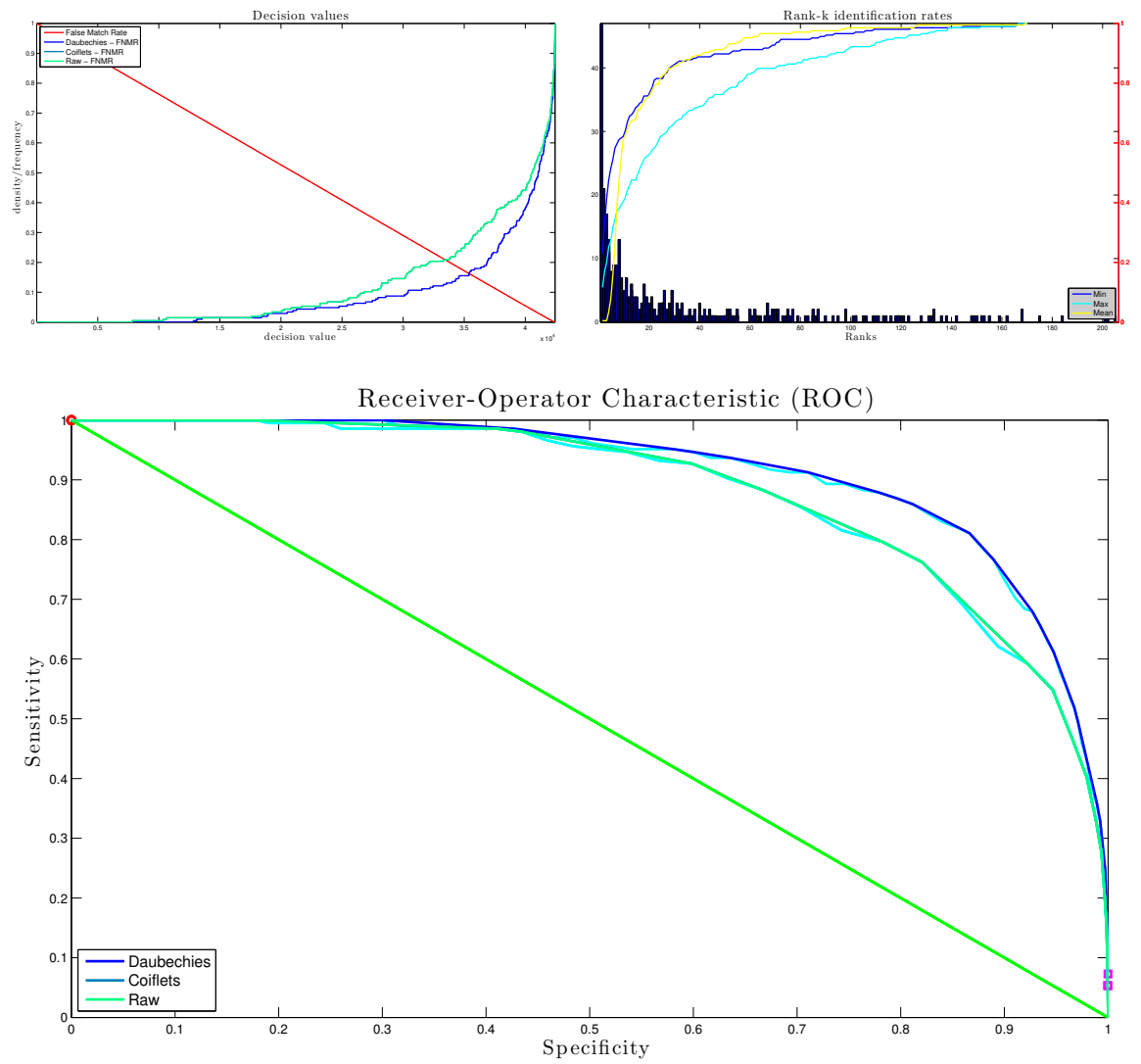


Figure 89: Influence of the wavelet denoising exhibited through the ROC curves of *Coiflets* and *Daubechies* wavelets.

a crucial part for the success of the biometric system (it can be related to the *template extraction* procedure). However, the segmentation algorithm is based on several choices for its parameters. Therefore, Figure 90 shows the influence of all segmentation parameters through the resulting ROC curves. Distinction is made between the choice of *positive* or *negative* difference in spectral flow, sizes of the analysis window of 0.7, 0.85 or 1 seconds with a number of partials between 0 and 2. We can see here that there is a whole span of different results depending on the settings of the segmentation procedure. First, the best results are obtained with the *positive* spectral flow for an analysis window of 0.85s and 2 partials connectivity used (these are the parameters used for the final analysis of results). Therefore, the corresponding curves are merged in the final figure. The widest disparities appears for the size of the analysis window. This parameter widely influence the final results, as it seems that a slightly smaller or wider size of window can dramatically change the accuracy. The two other parameter also seem to strongly influence the results. Therefore, a correct segmentation appears to be the most important step of the whole heart biometric system.

A.3.3 Beat Selection

After the segmentation procedure, a set of heart beats is obtained for each individual. However, this set may often contain some erroneous beats recordings, due to widely varying recordings. Therefore, we perform a selection of segmented beats in order to improve the homogeneity of the resulting set. This selection can be performed either on *energy deviation* (we select the heart beats which are almost of same total energy) and *shape deviation* (we select the set based on minimizing the inter-beats distance in energetic time series). Therefore, figure 91 exhibits the influence of each of these selection procedures based on their resulting ROC curves.

A.3.4 Time series comparison

Once all the heart beats are segmented, the temporal evolution of their features are compared thanks to different time series distances. We try to see the influence and differences between the use of DTW with different maximum warping windows and simply using the Euclidean distance on down-sampled time series.

Warping window

When using the DTW distance on time series, it is possible to constrain the maximum warping authorized between two series. This parameter allows to control the *reach* to which two points might be compared. Figure 92 shows the influence of the size of the warping window exhibited through the ROC curves of 2, 5, 10 and 20% of authorized warping reach.

Resampling

Another solution to compare the time series is to use the Euclidean distance. This distance is often claimed to lack the temporal flexibility required for precise matching. However, our large scale study on various datasets shows that the use of Euclidean distance on strongly down-sampled series often leads to better classification results than using the DTW. Figure 93 shows the influence of the size of the resampling factors

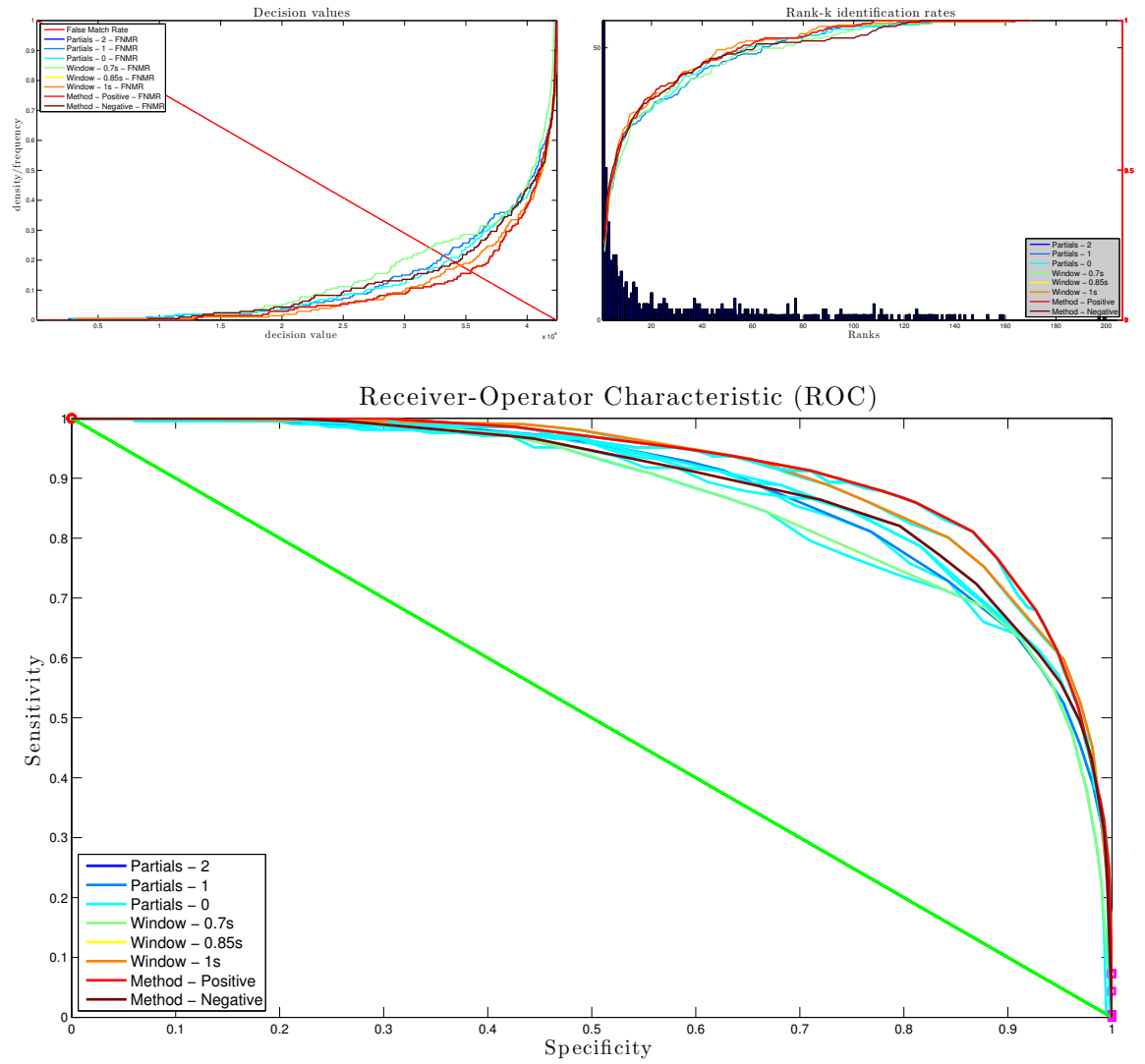


Figure 90: Influence of the segmentation parameters exhibited through the ROC curves of spectral differences (F+ or F-), window size and number of partials.

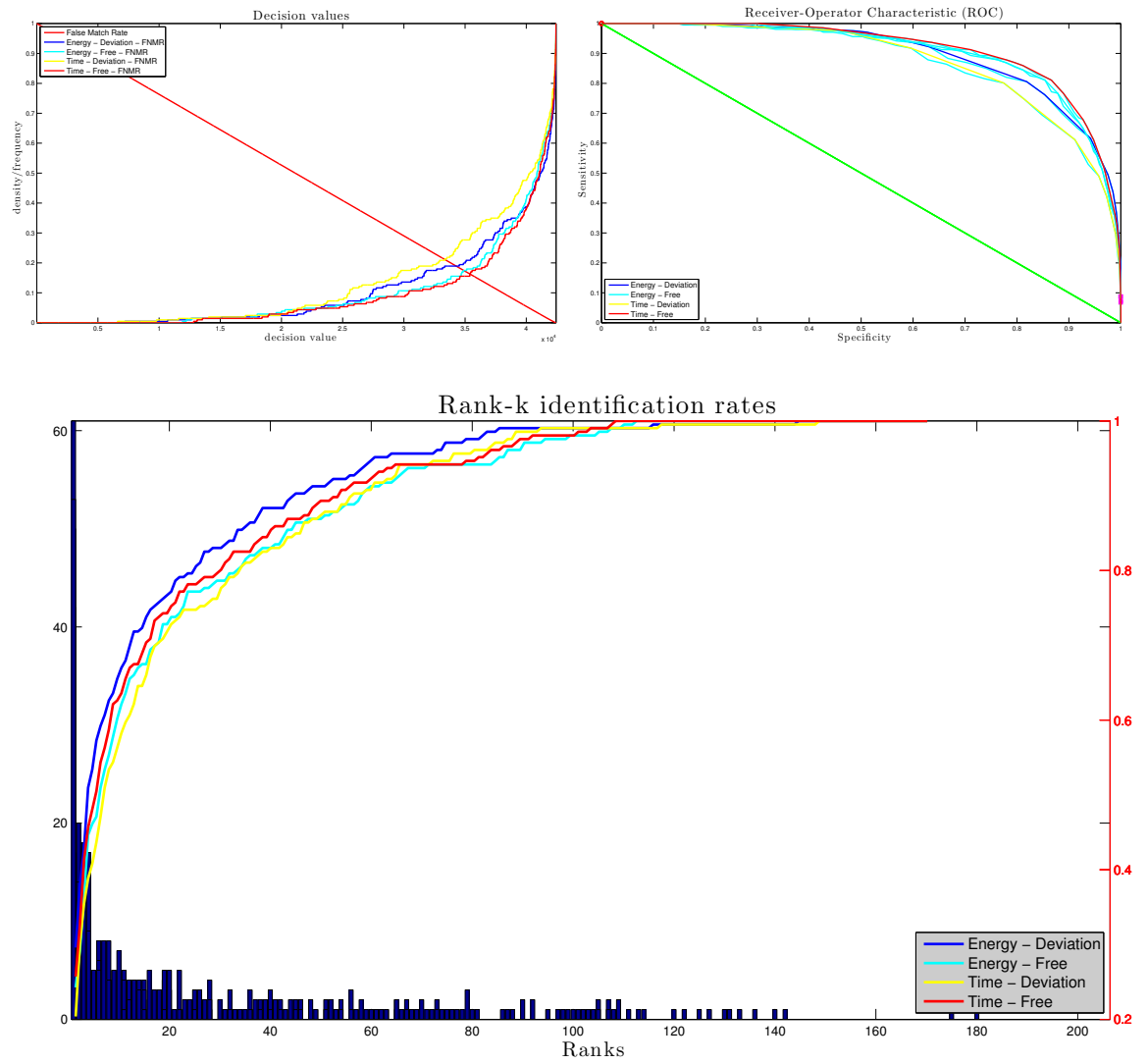


Figure 91: Influence of the beat selection exhibited through the ROC curves of *energy deviation* and *shape deviation*.

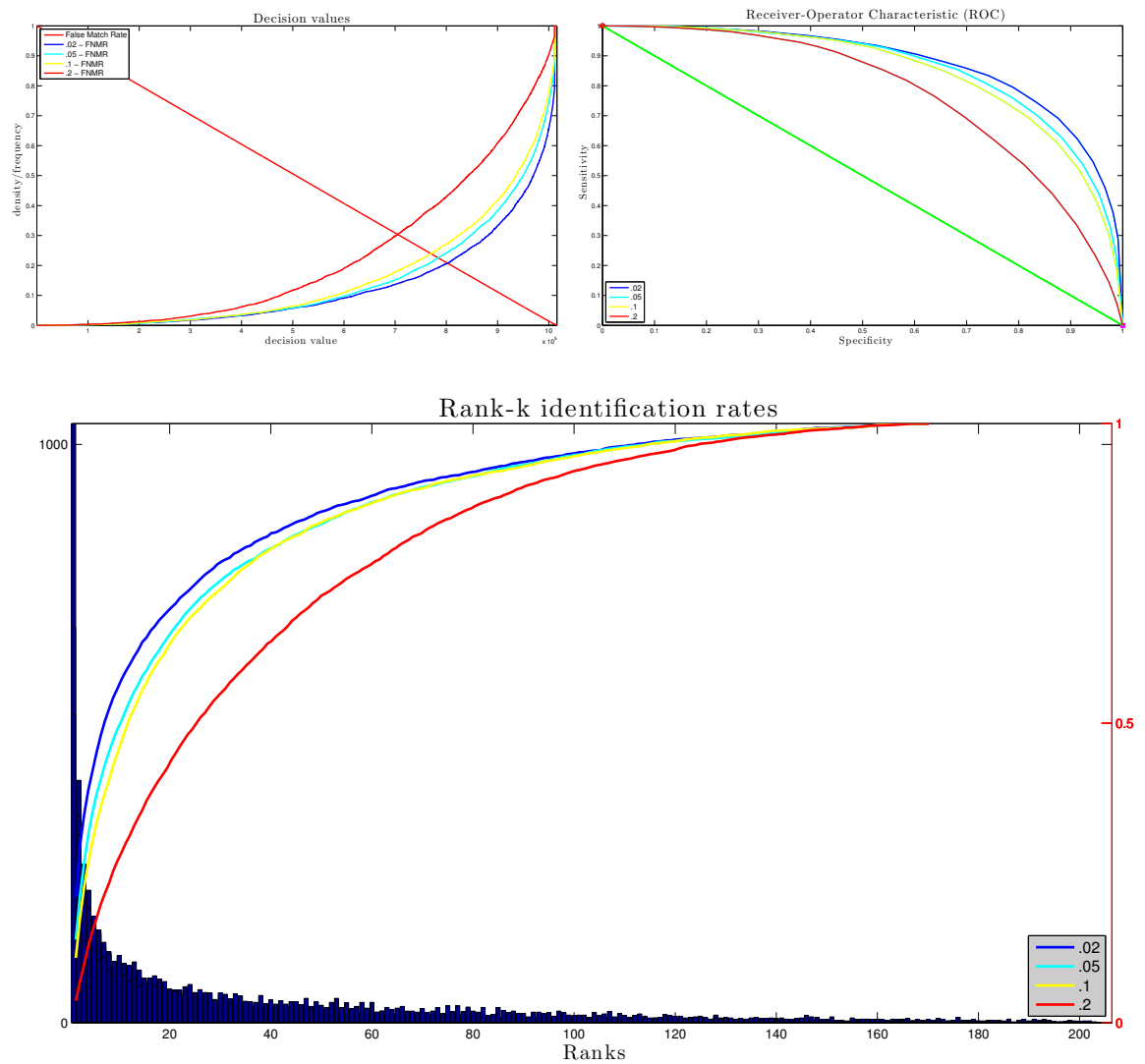


Figure 92: Influence of the size of the warping window for comparing the time series with the DTW distance measure exhibited through the ROC curves of 2, 5, 10 and 20% of authorized warping reach.

for comparing the time series exhibited through the ROC curves of 128, 64, 32, 16 and 8 resampling points.

A.3.5 *Decision influence*

We study here the influence of the decisions alternatives on the overall performances. This decision relies on the way the final similarity score is computed for one individual. First the type of *decision rule* decides how the score of multiple heart beats should be merged together. The number of heart beats to be merged is decided by the *size of testing set*. Each heartbeat is compared to the distribution in the database depending on the *size of training set*.

Type of decision rule

When a recording is input to the system, it provides several heartbeats to compare to the templates in the database. Therefore, we obtain a set of scores between each heartbeat and each class in the database. In order to make a final decision, we have therefore different possibilities of merging these scores. Figure 94 shows the influence of the type of *decision rule* exhibited through the ROC curves of the *mean*, *min* and *max* rules.

Size of testing set

The previous comparison made the assumption that all heart beats obtained from the segmentation and selection are used for the final score. However, as our selection procedure allows to rank each heart beat based on its *energy* and *shape* deviations, we try to see if we can obtain better results by constraining the maximum number of heart beats to use in the final scoring. Figure 95 shows the influence of the size of *testing set* exhibited through the ROC curves of different set cardinalities.

Size of training set

Based on the same idea as the previous comparison, we provide a comparison based on the size of the *training set*. Therefore, when computing the similarity scores, we reduce the size of the template database for each individual and use only a restricted set of templates. Figure 96 shows the influence of the size of *training set* exhibited through the ROC curves of different set cardinalities.

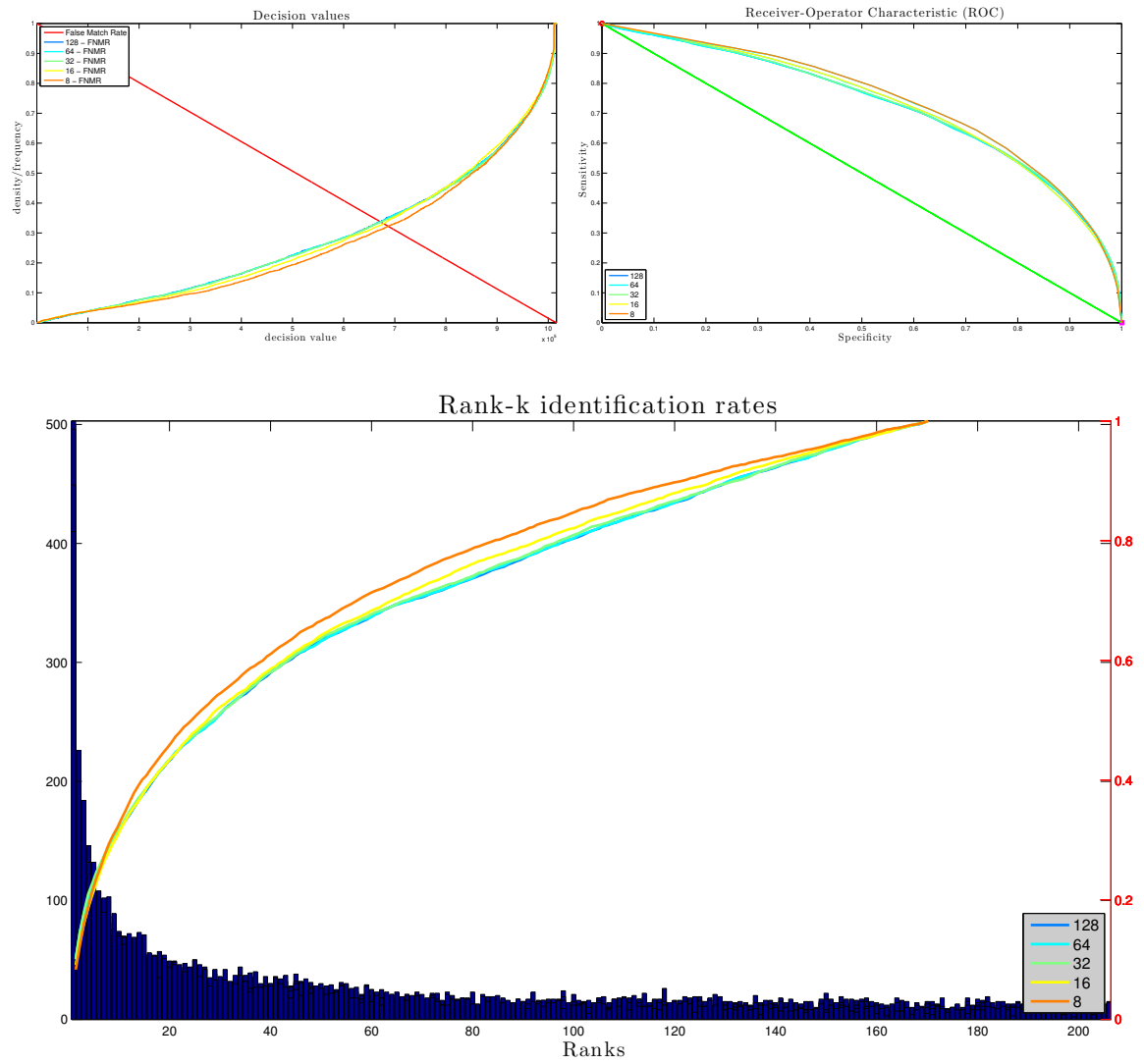


Figure 93: Influence of the size of the resampling factors for comparing the time series exhibited through the ROC curves of 128, 64, 32, 16 and 8 resampling points.

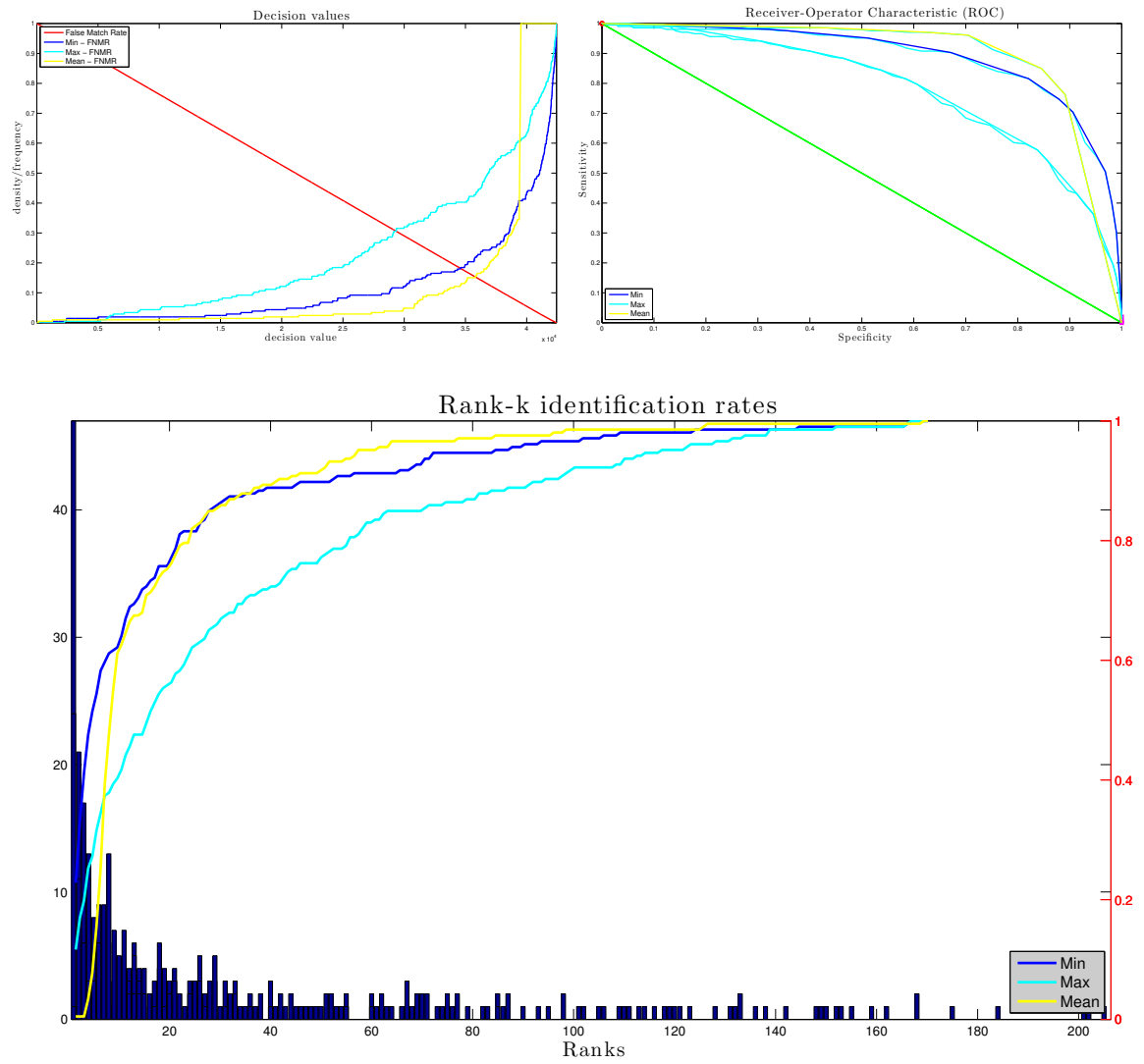


Figure 94: Influence of the type of *decision rule* exhibited through the ROC curves of the *mean*, *min* and *max* rules.

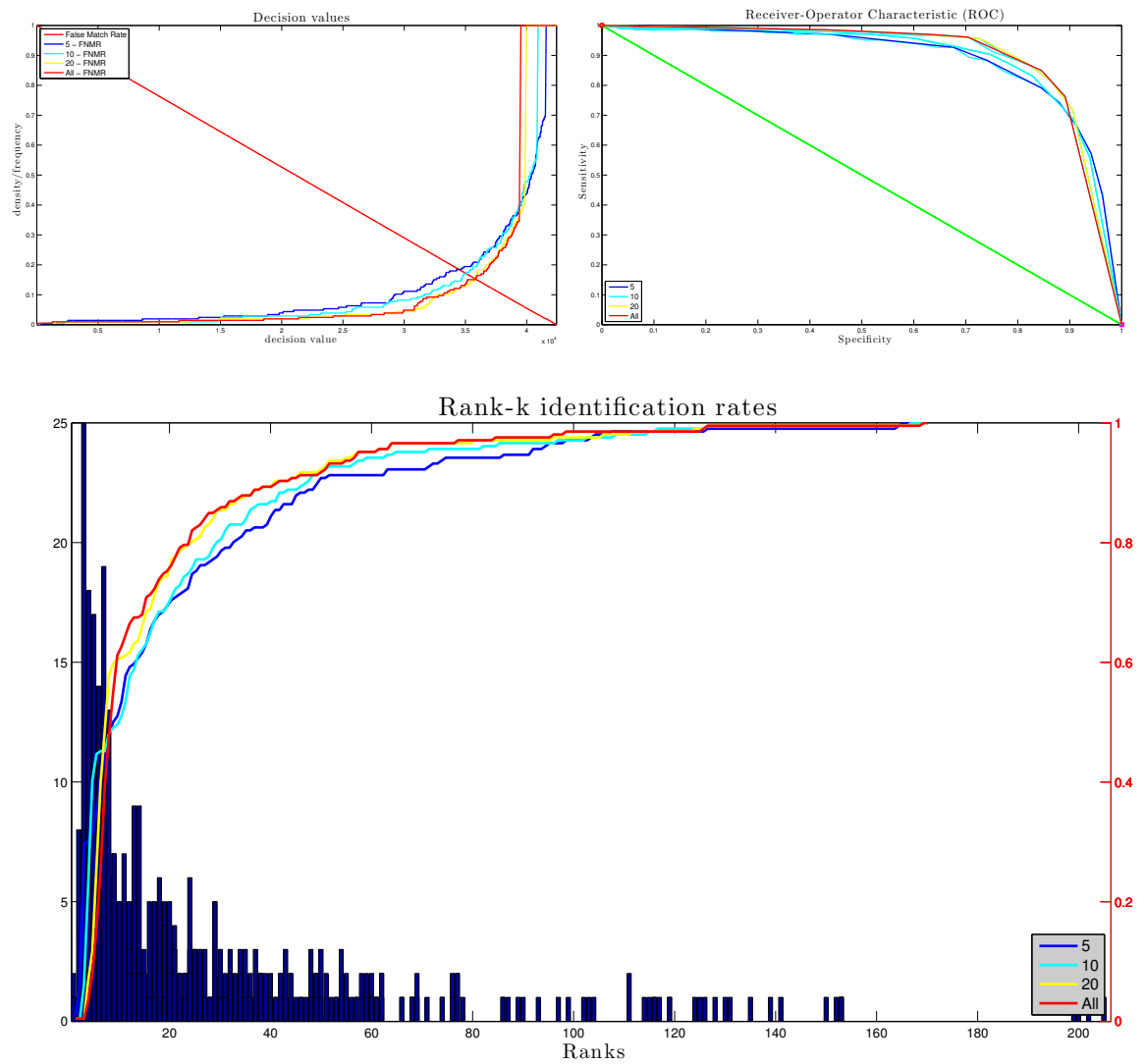


Figure 95: Influence of the size of *testing set* exhibited through the ROC curves of different set cardinalities.

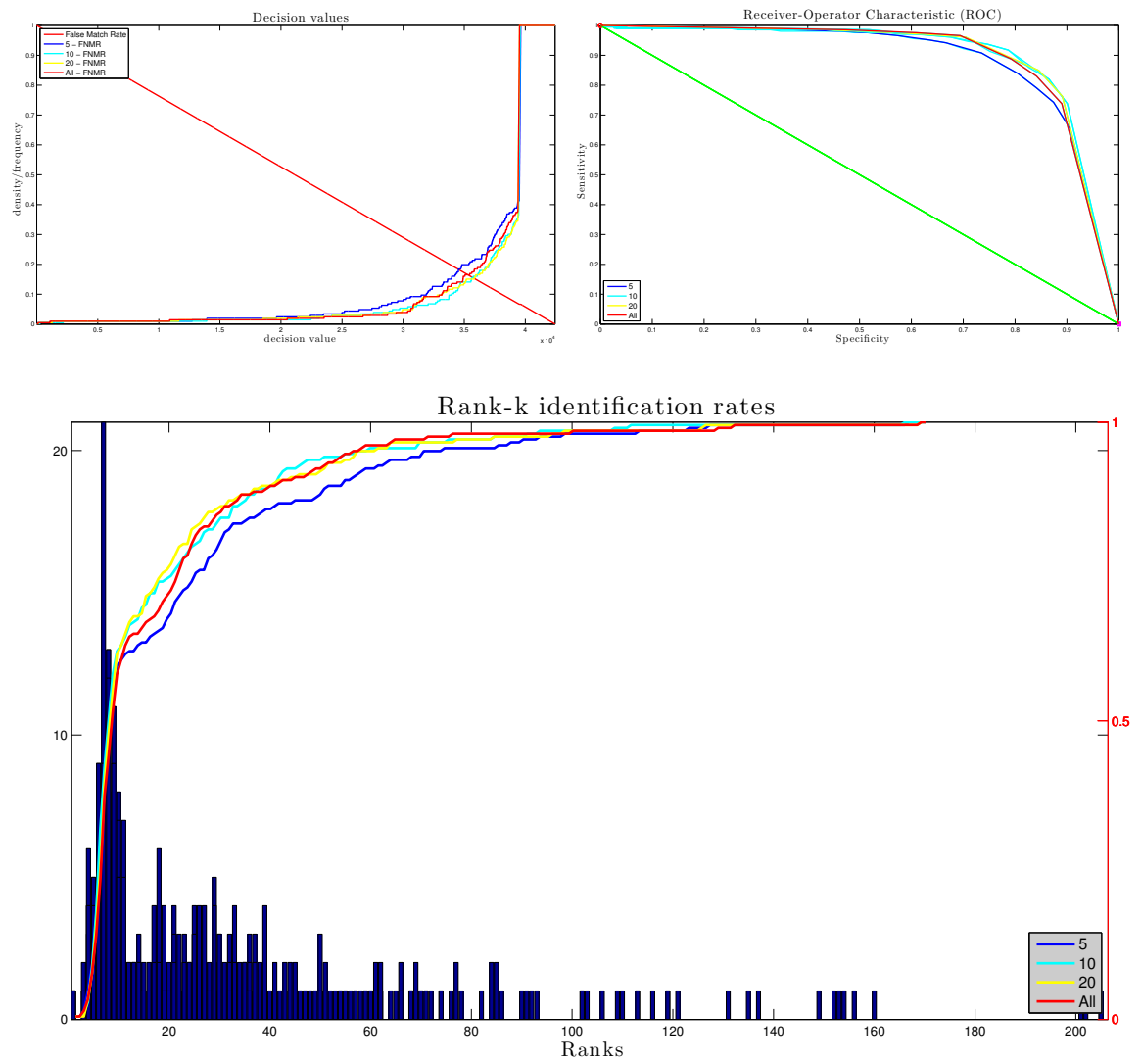


Figure 96: Influence of the size of *training set* exhibited through the ROC curves of different set cardinalities.

BIBLIOGRAPHY

- [1] ISO 9421-11. Ergonomic requirements for office work with visual display terminals - part 11: Guidance on usability. *International Organization of Standardization*, 1995. (Cited on page 79.)
- [2] J. Abonyi, B. Fell, S. Nemeth, and P. Arva. Fuzzy clustering based segmentation of time-series. In *Proceedings of the 5th International Symposium on Intelligent Data Analysis, IDA 2003, August 28-30*, pages 275–285, Berlin, Germany, 2003. Springer-Verlag, New York Inc. (Cited on page 30.)
- [3] M.F. Abu-Taleb and B. Mareschal. Water resources planning in the middle east: application of the PROMETHEE V multicriteria method. *European Journal of Operational Research*, 81(3):500–511, 1995. ISSN 0377-2217. (Cited on page 53.)
- [4] F. Agrafioti, F.M. Bui, and D. Hatzinakos. Medical biometrics: The perils of ignoring time dependency. In *Biometrics: Theory, Applications, and Systems, 2009. BTAS'09. IEEE 3rd International Conference on*, pages 1–6. IEEE, 2009. (Cited on page 165.)
- [5] R. Agrawal, C. Faloutsos, and A.N. Swami. Efficient Similarity Search In Sequence Databases. In *Proceedings of the 4th International Conference on Foundations of Data Organization and Algorithms*, pages 69–84, Chicago, Illinois, USA, 1993. Springer. (Cited on pages 25 and 36.)
- [6] Rakesh Agrawal, King-Ip Lin, Harpreet S. Sawhney, and Kyuseok Shim. Fast similarity search in the presence of noise, scaling, and translation in time-series databases. In *VLDB '95: Proceedings of the 21th International Conference on Very Large Data Bases*, pages 490–501, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1-55860-379-4. (Cited on pages 26 and 34.)
- [7] N.K. Ahmed, A.F. Atiya, N. El Gayar, H. El-Shishiny, and E. Giza. An empirical comparison of machine learning models for time series forecasting. *Econometric Reviews*, 29(5):594–621, 2009. (Cited on page 31.)
- [8] T. Ahmed, B. Oreshkin, and M. Coates. Machine learning approaches to network anomaly detection. In *Proceedings of the 2nd USENIX workshop on Tackling computer systems problems with machine learning techniques*, pages 1–6, Cambridge, MA, USA, 2007. USENIX Association. (Cited on page 32.)
- [9] F. Alimoglu and E. Alpaydin. Methods of combining multiple classifiers based on different representations for pen-based handwritten digit recognition. In *Proceedings of the Fifth Turkish Artificial Intelligence and Artificial Neural Networks Symposium*. Citeseer, 1996. (Cited on pages 130 and 223.)
- [10] J. An, H. Chen, K. Furuse, N. Ohbo, and E. Keogh. Grid-based indexing for large time series databases. *Intelligent Data Engineering and Automated Learning, Lecture Notes in Computer Science*, 1983(1):614–621, 2003. (Cited on page 37.)
- [11] C.M. Antunes and A.L. Oliveira. Temporal data mining: An overview. In *KDD Workshop on Temporal Data Mining*, pages 1–13, San Francisco, CA, USA, 2001. (Cited on pages 22, 39, and 42.)

- [12] T. Argyros and C. Ermopoulos. Efficient subsequence matching in time series databases under time and amplitude transformations. In *3rd IEEE International Conference on Data Mining*, pages 481–484, 2003. (Cited on page 34.)
- [13] I. Assent, R. Krieger, F. Afschari, and T. Seidl. The TS-tree: efficient time series search and retrieval. In *Proceedings of the 11th International Conference on Extending Database Technology*, pages 25–29, 2008. (Cited on page 46.)
- [14] I. Assent, M. Wichterich, R. Krieger, H. Kremer, and T. Seidl. Anticipatory DTW for efficient similarity search in time series databases. *Proceedings of the VLDB Endowment*, 2(1):826–837, 2009. (Cited on page 26.)
- [15] J. Aßfalg, H.P. Kriegel, P. Kroger, P. Kunath, A. Pryakhin, and M. Renz. Similarity search on time series based on threshold queries. In *Advances in database technology: EDBT 2006: 10th International Conference on Extending Database Technology, March 26–31*, volume 3896, page 276, Munich, Germany, 2006. Springer-Verlag New York Inc. (Cited on page 42.)
- [16] J. Aßfalg, H.P. Kriegel, P. Kröger, P. Kunath, A. Pryakhin, and M. Renz. Similarity search in multimedia time series data using amplitude-level features. In *Proceedings of the 14th international conference on Advances in multimedia modeling*, pages 123–133. Springer-Verlag, 2008. (Cited on page 36.)
- [17] American Standards Association. *American Standard Acoustical Terminology*. ASA, New York, 1960. (Cited on pages 7 and 9.)
- [18] A. Bagnall and G. Janacek. Clustering time series with clipped data. *Machine Learning*, 58(2):151–178, 2005. (Cited on page 27.)
- [19] A. Bagnall, C.A. Ratanamahatana, E. Keogh, S. Lonardi, and G. Janacek. A bit level representation for time series data mining with shape based similarity. *Data mining and knowledge discovery*, 13(1):11–40, 2006. (Cited on page 37.)
- [20] AJ Bagnall, G. Janacek, B. De la Iglesia, and M. Zhang. Clustering time series from mixture polynomial models with discretised data. In *Proceedings of the 2nd Australasian Data Mining Workshop*, pages 105–120, 2003. (Cited on page 37.)
- [21] J. Bai and S. Ng. Forecasting economic time series using targeted predictors. *Journal of Econometrics*, 146(2):304–317, 2008. (Cited on page 31.)
- [22] BR Bakshi and G. Stephanopoulos. Representation of process trends–IV. Induction of real-time patterns from operating data for diagnosis and supervisory control. *Computers & Chemical Engineering*, 18(4):303–332, 1994. (Cited on pages 22 and 28.)
- [23] B.R. Bakshi and G. Stephanopoulos. Reasoning in time: Modeling, analysis, and pattern recognition of temporal process trends. *Advances in Chemical Engineering*, 22:485–548, 1995. (Cited on page 37.)
- [24] G Ballet, R Borghesi, P Hoffmann, and F Levy. Studio online 3.0 : An internet "killer application" for remote access to ircam sounds and processing tools. In *Actes des Journees Informatique Musicale*, Paris, France, 1999. (Cited on page 66.)
- [25] JP Bandera, R. Marfil, A. Bandera, JA Rodríguez, L. Molina-Tanco, and F. Sandoval. Fast gesture recognition based on a two-level representation. *Pattern Recognition Letters*, 30(13):1181–1189, 2009. (Cited on page 37.)

- [26] P. Barone, M.F. Carfora, and R. March. Segmentation, Classification and Denoising of a Time Series Field by a Variational Method. *Journal of Mathematical Imaging and Vision*, 34(2):152–164, 2009. (Cited on page 48.)
- [27] G.A. Barreto. Time Series Prediction with the Self-Organizing Map: A Review. *Perspectives of neural-symbolic integration*, 77(1):135–158, 2007. (Cited on page 31.)
- [28] L. Barrington, A. Chan, D. Turnbull, and G. Lanckriet. Audio information retrieval using semantic similarity. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2007)*, volume 2, pages 722–725, 2007. (Cited on pages 71 and 167.)
- [29] I. Bartolini, P. Ciaccia, and M. Patella. Warp: Accurate retrieval of shapes using phase of fourier descriptors and time warping distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(1):142–147, 2005. (Cited on page 42.)
- [30] R. Bayer and EM McCreight. Organization and maintenance of large ordered indexes. *Acta informatica*, 1(3):173–189, 1972. (Cited on page 45.)
- [31] N. Beckmann, H.P. Kriegel, R. Schneider, and B. Seeger. The R*-tree: an efficient and robust access method for points and rectangles. *ACM SIGMOD Record*, 19(2):322–331, 1990. (Cited on pages 25 and 45.)
- [32] N. Belacel. Multicriteria assignment method PROAFTN: Methodology and medical application. *European Journal of Operational Research*, 125(1):175–183, 2000. ISSN 0377-2217. (Cited on pages 49 and 53.)
- [33] N. Belacel and M.R. Boulassel. Multicriteria fuzzy assignment method: a useful tool to assist medical diagnosis. *Artificial intelligence in medicine*, 21(1-3):201–207, 2001. ISSN 0933-3657. (Cited on page 53.)
- [34] N. Belacel, P. Vincke, JM Scheiff, and MR Boulassel. Acute leukemia diagnosis aid using multicriteria fuzzy assignment methodology. *Computer methods and programs in biomedicine*, 64(2):145, 2001. ISSN 0169-2607. (Cited on page 53.)
- [35] S. Berchtold, D.A. Keim, and H.P. Kriegel. The X-tree: An index structure for high-dimensional data. *Readings in multimedia computing and networking*, 4(1):451–463, 2002. (Cited on page 45.)
- [36] F. Beritelli and S. Serrano. Biometric identification based on frequency analysis of cardiac sounds. *Information Forensics and Security, IEEE Transactions on*, 2(3):596–604, 2007. (Cited on page 238.)
- [37] F. Beritelli and A. Spadaccini. Human identity verification based on mel frequency analysis of digital heart sounds. In *Digital Signal Processing, 2009 16th International Conference on*, pages 1–5. IEEE, 2009. (Cited on page 238.)
- [38] F. Beritelli and A. Spadaccini. An improved biometric identification system based on heart sounds and gaussian mixture models. In *Biometric Measurements and Systems for Security and Medical Applications (BIOMS), 2010 IEEE Workshop on*, pages 31–35. IEEE, 2010. (Cited on pages 147, 153, 161, and 238.)
- [39] P. Berkhin. A survey of clustering data mining techniques. *Grouping Multidimensional Data*, pages 25–71, 2006. (Cited on page 27.)

- [40] H. Berlioz. *Traité d'instrumentation et d'orchestration*. Gregg International Publishers, 1855. (Cited on pages 2, 5, and 8.)
- [41] D. Berndt and J. Clifford. Using dynamic time warping to find patterns in time series. In *AAAI-94 workshop on knowledge discovery in databases*, pages 229–248, 1994. (Cited on page 39.)
- [42] T. Bernecker, M. Houle, H.P. Kriegel, P. Kroger, M. Renz, E. Schubert, and A. Zimek. Quality of similarity rankings in time series. *Advances in Spatial and Temporal Databases*, pages 422–440, 2011. (Cited on pages 130 and 224.)
- [43] S. Berretti, A. Del Bimbo, and P. Pala. Retrieval by shape similarity with perceptual distance and effective indexing. *IEEE Transactions on multimedia*, 2(4): 225–239, 2000. (Cited on page 48.)
- [44] R. Bhargava, H. Kargupta, and M. Powers. Energy consumption in data analysis for on-board and distributed applications. In *Proceedings of the ICML*, volume 3, 2003. (Cited on page 47.)
- [45] V. Bhaskar, S.K. Gupta, and A.K. Ray. Applications of multiobjective optimization in chemical engineering. *Reviews in Chemical Engineering*, 16(1):1–54, 2000. ISSN 0167-8299. (Cited on pages 49 and 53.)
- [46] M. Bicego, V. Murino, and M. Figueiredo. Similarity-based clustering of sequences using hidden Markov models. *Lecture Notes in Computer Science*, 2743:95–104, 2003. (Cited on page 42.)
- [47] L. Biel, O. Pettersson, L. Philipson, and P. Wide. Ecg analysis: a new approach in human identification. *Instrumentation and Measurement, IEEE Transactions on*, 50(3):808–812, 2001. (Cited on page 147.)
- [48] V.L. Blankers, C. Heuvel, K.Y. Franke, and L.G. Vuurpijl. Icdar 2009 signature verification competition. In *Document Analysis and Recognition, 2009. ICDAR'09. 10th International Conference on*, pages 1403–1407. IEEE, 2009. (Cited on page 163.)
- [49] B. Blankertz, K.R. Muller, G. Curio, T.M. Vaughan, G. Schalk, J.R. Wolpaw, A. Schlogl, C. Neuper, G. Pfurtscheller, T. Hinterberger, et al. The bci competition 2003: progress and perspectives in detection and discrimination of eeg single trials. *Biomedical Engineering, IEEE Transactions on*, 51(6):1044–1051, 2004. (Cited on pages 130, 218, and 219.)
- [50] B. Blankertz, G. Dornhege, M. Krauledat, K.R. Muller, and G. Curio. The non-invasive berlin brain-computer interface: Fast acquisition of effective performance in untrained subjects. *NeuroImage*, 37(2):539–550, 2007. (Cited on pages 130 and 220.)
- [51] C. Bohm, S. Berchtold, and D.A. Keim. Searching in high-dimensional spaces: Index structures for improving the performance of multimedia databases. *ACM Computing Surveys*, 33(3):322–373, 2001. (Cited on page 45.)
- [52] B. Bollobas, G. Das, D. Gunopulos, and H. Mannila. Time-series similarity problems and well-separated geometric sets. In *Proceedings of the 13th symposium on computational geometry*, pages 454–456, 1997. (Cited on page 41.)

- [53] M. Botta, A. Giordana, and L. Saitta. Learning fuzzy concept definitions. In *Fuzzy Systems, 1993., Second IEEE International Conference on*, pages 18–22. IEEE, 1993. (Cited on pages 130 and 216.)
- [54] R. Bousseljot, D. Kreiseler, and A. Schnabel. Nutzung der ekg-signaldatenbank cardiodat der ptb uber das internet. *Biomedizinische Technik/Biomedical Engineering*, 40(s1):317–318, 1995. (Cited on page 231.)
- [55] K.W. Bowyer, K. Chang, and P. Flynn. A survey of approaches and challenges in 3d and multi-modal 3d+ 2d face recognition. *Computer Vision and Image Understanding*, 101(1):1–15, 2006. (Cited on page 163.)
- [56] G.E.P. Box, G.M. Jenkins, and G.C. Reinsel. *Time series analysis: forecasting and control*. Holden-day San Francisco, 1976. (Cited on page 30.)
- [57] A. Brancucci and P. San Martini. Laterality in the perception of temporal cues of musical timbre. *Neuropsychologia*, 37(13):1445–1451, 1999. (Cited on page 10.)
- [58] AS Bregman. Timbre, orchestration, dissonance, et organisation auditive. *Le timbre, métaphore pour la composition*, pages 204–215, 1991. (Cited on pages 4 and 8.)
- [59] P.J. Brockwell and R.A. Davis. *Introduction to time series and forecasting*. Springer Verlag, 2002. (Cited on page 30.)
- [60] P.J. Brockwell and R.A. Davis. *Time series: theory and methods*. Springer Verlag, 2009. ISBN 1441903194. (Cited on page 30.)
- [61] R.A. Brown and R. Frayne. A fast discrete s-transform for biomedical signal processing. In *Engineering in Medicine and Biology Society, 2008. EMBS 2008. 30th Annual International Conference of the IEEE*, pages 2586–2589. IEEE, 2008. (Cited on pages 151 and 239.)
- [62] M. Bsoul, H. Minn, M. Nourani, G. Gupta, and L. Tamil. Real-time sleep quality assessment using single-lead ecg and multi-stage svm classifier. In *Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conference*, volume 2010, page 1178, 2010. (Cited on pages 130 and 234.)
- [63] J. Buhler and M. Tompa. Finding motifs using random projections. *Journal of computational biology*, 9(2):225–242, 2002. ISSN 1066-5277. (Cited on page 33.)
- [64] H.S. Burkom, S.P. Murphy, and G. Shmueli. Automated time series forecasting for biosurveillance. *Statistics in Medicine*, 26(22):4202–4218, 2007. (Cited on pages 22 and 31.)
- [65] R. Cai, L. Lu, and H.J. Zhang. Using structure patterns of temporal and spectral feature in audio similarity measure. In *Proceedings of the 11th ACM international conference on Multimedia*, pages 219–222, 2003. (Cited on pages 71 and 167.)
- [66] Y. Cai and R. Ng. Indexing spatio-temporal trajectories with Chebyshev polynomials. In *Proceedings of the 2004 ACM SIGMOD international conference on Management of data*, pages 599–610, Paris, France, 2004. ACM. (Cited on page 36.)

- [67] L.M. Camarinha-Matos, L.S. Lopes, and J. Barata. Integration and learning in supervision of flexible assembly systems. *Robotics and Automation, IEEE Transactions on*, 12(2):202–219, 1996. (Cited on pages 130 and 232.)
- [68] P. Cano and M. Koppenberger. Automatic sound annotation. In *Proceedings of the 14th IEEE Workshop on Machine Learning for Signal Processing*, pages 391–400. IEEE, 2004. (Cited on pages 167 and 173.)
- [69] P. Cano, M. Koppenberger, S. Le Groux, J. Ricard, P. Herrera, and N. Wack. Nearest-neighbor generic sound classification with a wordnet-based taxonomy. In *Proceedings of the 116th AES Convention*, Berlin, Germany, 2004. Citeseer. (Cited on pages 71 and 167.)
- [70] L. Cao and F. Tay. Feature selection for support vector machines in financial time series forecasting. *Intelligent Data Engineering and Automated Learning. Lecture Notes in Computer Science*, 1983:41–65, 2009. (Cited on page 31.)
- [71] R. Cappelli, M. Ferrara, and D. Maltoni. Fingerprint indexing based on minutia cylinder-code. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(5):1051–1057, 2011. (Cited on page 163.)
- [72] G. Carpentier. *Approche computationnelle de l'orchestration musicale-Optimisation multicritère sous contraintes de combinaisons instrumentales dans de grandes banques de sons*. PhD thesis, Université Pierre et Marie Curie-Paris VI, 2008. (Cited on pages 185, 190, and 192.)
- [73] G. Carpentier, D. Tardieu, G. Assayag, X. Rodet, and E. Saint-James. Imitative and generative orchestrations using pre-analyzed sound databases. In *Proceedings of Sound and Music Computing Conference (SMC)(Marseille, France)*, pages 115–122, 2006. (Cited on page 185.)
- [74] A. Casella, V. Mortari, and P. Petit. *La Technique de l'orchestre contemporain*. Ricordi, 1958. (Cited on page 2.)
- [75] M. Casey. General sound classification and similarity in mpeg-7. *Organised Sound*, 6(02):153–164, 2001. (Cited on page 167.)
- [76] M.A. Casey. Reduced-rank spectra and minimum-entropy priors as consistent and reliable cues for generalized sound recognition. In *Workshop for Consistent & Reliable Acoustic Cues*, 2001. (Cited on page 167.)
- [77] M.A. Casey. Sound classification and similarity. *Introduction to MPEG-7: Multimedia Content Description Interface*, pages 309–317, 2002. (Cited on page 167.)
- [78] M.A. Casey. Acoustic lexemes for organizing internet audio. *Contemporary Music Review*, 24(6):489–508, 2005. (Cited on page 71.)
- [79] M.A. Casey, R. Veltkamp, M. Goto, M. Leman, C. Rhodes, and M. Slaney. Content-based music information retrieval: current directions and future challenges. *Proceedings of the IEEE*, 96(4):668–696, 2008. (Cited on page 70.)
- [80] M.J. Castro-Bleda, S. Espana, J. Gorbe, F. Zamora, D. Llorens, A. Marzal, F. Prat, and J.M. Vilar. Improving a dtw-based recognition engine for on-line handwritten characters by using mlps. In *10th International Conference on Document Analysis and Recognition*, pages 1260–1264, 2009. (Cited on pages 130 and 229.)

- [81] K. Chakrabarti and S. Mehrotra. The hybrid tree: an index structure for high dimensional feature spaces. In *Data Engineering, 1999. Proceedings., 15th International Conference on*, pages 440–447, 1999. (Cited on page 45.)
- [82] F.K.P. Chan, A.W.C. Fu, and C. Yu. Haar wavelets for efficient similarity search of time-series: with and without time warping. *IEEE Transactions on knowledge and data engineering*, 15(3):686–705, 2003. (Cited on page 36.)
- [83] K Chan and AW Fu. Efficient time series matching by wavelets. In *Proceedings of the 15th IEEE International conference on data engineering*, pages 126 – 133, Sydney, Australia, 1999. (Cited on pages 25 and 36.)
- [84] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Computing Surveys*, 41(3):15, 2009. (Cited on page 32.)
- [85] V. Chankong and Y.Y. Haimes. *Multiobjective decision making: theory and methodology*. North-Holland New York, 1983. ISBN 0444007105. (Cited on page 49.)
- [86] JC Chappelier and A. Grumbach. A Kohonen map for temporal sequences. In *Proceedings of the Conference on Neural Networks and Their Applications.*, pages 104–110, 1996. (Cited on page 27.)
- [87] L. Chen and R. Ng. On the marriage of Lp-norms and edit distance. In *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30*, pages 792–803. VLDB Endowment, 2004. (Cited on page 41.)
- [88] L. Chen, M.T. Ozsu, and V. Oria. Robust and fast similarity search for moving object trajectories. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pages 491–502, Baltimore, Maryland, USA, 2005. ACM. (Cited on page 41.)
- [89] Q. Chen, L. Chen, X. Lian, Y. Liu, and J.X. Yu. Indexable PLA for efficient similarity search. In *Proceedings of the 33rd international conference on Very large data bases*, pages 435–446. VLDB Endowment, 2007. (Cited on page 37.)
- [90] X. Chen and Y. Zhan. Multi-scale anomaly detection algorithm based on infrequent pattern of time series. *Journal of Computational and Applied Mathematics*, 214(1):227–237, 2008. (Cited on page 32.)
- [91] X. Chen, S. Kwong, and M. Li. A compression algorithm for DNA sequences and its applications in genome comparison. In *Proceedings of the fourth annual international conference on Computational molecular biology*, page 107, 2000. (Cited on page 43.)
- [92] Y. Chen, MA Nascimento, B.C. Ooi, and AKH Tung. Spade: On shape-based pattern detection in streaming time series. In *IEEE 23rd International Conference on Data Engineering, 2007*, pages 786–795, 2007. (Cited on page 41.)
- [93] V. Chhieng and R. Wong. Adaptive distance measurement for time series databases. *Lecture Notes in Computer Science*, 4443:598–610, 2010. (Cited on page 41.)
- [94] B. Chiu, E. Keogh, and S. Lonardi. Probabilistic discovery of time series motifs. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 493–498, Washington, D.C, USA, 2003. ACM. (Cited on page 33.)

- [95] M. Chuah and F. Fu. ECG anomaly detection via time series analysis. In *Frontiers of High Performance Computing and Networking ISPA 2007 Workshops*, pages 123–135. Springer, 2007. (Cited on pages 31, 32, and 46.)
- [96] W.J. Conover. Practical nonparametric statistics, 1980. (Cited on pages 131 and 170.)
- [97] D. Coomans and DL Massart. Alternative k -nearest neighbour rules in supervised pattern recognition: Part 1. k -nearest neighbour classification by using alternative voting rules. *Analytica Chimica Acta*, 136:15–27, 1982. (Cited on page 126.)
- [98] M. Corduas and D. Piccolo. Time series clustering and classification by the autoregressive metric. *Computational statistics & data analysis*, 52(4):1860–1872, 2008. (Cited on page 27.)
- [99] G. Cormode, S. Muthukrishnan, and W. Zhuang. Conquering the divide: Continuous clustering of distributed data streams. In *IEEE 23rd International Conference on Data Engineering*, 2007, pages 1036–1045, 2007. (Cited on page 27.)
- [100] C. Costa Santos, J. Bernardes, PMB Vitanyi, and L. Antunes. Clustering fetal heart rate tracings by compression. In *19th International Symposium on Computer-Based Medical Systems*, pages 685–690, 2006. (Cited on page 43.)
- [101] D. Coyle, G. Prasad, and T.M. McGinnity. Extracting features for a brain-computer interface by self-organising fuzzy neural network-based time series prediction. In *Engineering in Medicine and Biology Society, 2004. IEMBS'04. 26th Annual International Conference of the IEEE*, volume 2, pages 4371–4374. IEEE, 2004. (Cited on pages 130 and 221.)
- [102] G. Das, D. Gunopulos, and H. Mannila. Finding similar time series. In *Principles of data mining and knowledge discovery: First European Symposium, PKDD'97, June 24-27*, volume 1263, pages 88–100, Trondheim, Norway, 1997. Springer Verlag. (Cited on page 41.)
- [103] S. Davis, A. Jacobson, D. Suszcynsky, M. Cai, and D. Eads. Forte time series dataset, 2005. URL <http://nis-www.lanl.gov/eads/datasets/forte>. (Cited on page 224.)
- [104] M. Degli Esposti, C. Farinelli, and G. Menconi. Sequence distance via parsing complexity: Heartbeat signals. *Chaos, Solitons & Fractals*, 39(3):991–999, 2009. (Cited on page 43.)
- [105] J. Demsar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research*, 7:1–30, 2006. (Cited on pages 129, 131, 134, and 204.)
- [106] K. Deng, AW Moore, and MC Nechyba. Learning to recognize time series: combining ARMA models with memory-based learning. In *Proceedings of the IEEE International Symposium on Computational Intelligence in Robotics and Automation*, 1997. CIRA'97, pages 246–251, 1997. (Cited on page 29.)
- [107] A. Denton. Kernel-density-based clustering of time series subsequences using a continuous random-walk noise model. In *Proceedings of the fifth IEEE International Conference on Data Mining*, pages 122–129, 2005. (Cited on page 28.)

- [108] D.B. Dias, R.C.B. Madeo, T. Rocha, H.H. Biscaro, and S.M. Peres. Hand movement recognition for brazilian sign language: A study using distance-based neural networks. In *Neural Networks, 2009. IJCNN 2009. International Joint Conference on*, pages 697–704. IEEE, 2009. (Cited on page 227.)
- [109] H. Ding, G. Trajcevski, P. Scheuermann, X. Wang, and E. Keogh. Querying and mining of time series data: experimental comparison of representations and distance measures. *Proceedings of the VLDB Endowment*, 1(2):1542–1552, 2008. (Cited on pages 36, 38, 39, 42, 43, 48, and 126.)
- [110] P. Domingos and G. Hulten. Mining high-speed data streams. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 71–80. ACM, 2000. (Cited on page 46.)
- [111] G. Dong, J. Han, L.V.S. Lakshmanan, J. Pei, H. Wang, and P.S. Yu. Online mining of changes from data streams: Research problems and preliminary results. In *Proceedings of the 2003 ACM SIGMOD Workshop on Management and Processing of Data Streams*, San Diego, CA, 2003. (Cited on page 48.)
- [112] B. Dorizzi, R. Cappelli, M. Ferrara, D. Maio, D. Maltoni, N. Houmani, S. Garcia-Salicetti, and A. Mayoue. Fingerprint and on-line signature verification competitions at icb 2009. *Advances in Biometrics*, pages 725–732, 2009. (Cited on page 163.)
- [113] G. Dornhege, B. Blankertz, G. Curio, and K.R. Muller. Boosting bit rates in noninvasive eeg single-trial classifications by feature combination and multiclass paradigms. *Biomedical Engineering, IEEE Transactions on*, 51(6):993–1002, 2004. (Cited on page 219.)
- [114] J.S. Downie. The scientific evaluation of music information retrieval systems: Foundations and future. *Computer Music Journal*, 28(2):12–23, 2004. (Cited on pages 14, 56, 71, 72, 84, and 131.)
- [115] J.S. Downie and S.J. Cunningham. Toward a theory of music information retrieval queries: System design implications. In *Proceedings: Third International Conference on Music Information Retrieval. ISMIR*, volume 2002, pages 13–17, 2002. (Cited on pages 56, 84, and 90.)
- [116] H.W. Draper, C.J. Pfeffer, F.W. Stallmann, D. Littmann, and H.V. Pipberger. The corrected orthogonal electrocardiogram and vectorcardiogram in 510 normal men (frank lead system). *Circulation*, 30(6):853–864, 1964. (Cited on page 164.)
- [117] F.Y. Edgeworth. Mathematical psychics. *History of Economic Thought Books*, 1881. (Cited on page 49.)
- [118] J. Eggermont, J.N. Kok, and W.A. Kusters. Genetic programming for data classification: Partitioning the search space. In *Proceedings of the 2004 ACM symposium on Applied computing*, pages 1001–1005. ACM, 2004. (Cited on pages 130 and 226.)
- [119] M. Ehrgott. *Multicriteria optimization*, volume 491. Springer Verlag, 2005. (Cited on page 49.)
- [120] M. Ehrgott and X. Gandibleux. A survey and annotated bibliography of multiobjective combinatorial optimization. *Or Spectrum*, 22(4):425–460, 2000. (Cited on page 53.)

- [121] N. El-Bendary, H. Al-Qaheri, HM Zawbaa, M. Hamed, A.E. Hassanien, Q. Zhao, and A. Abraham. Hsas: Heart sound authentication system. In *Nature and Biologically Inspired Computing (NaBIC), 2010 Second World Congress on*, pages 351–356. IEEE, 2010. (Cited on pages 147 and 238.)
- [122] G. Ellis and A. Dix. An explorative analysis of user evaluation studies in information visualisation. In *Proceedings of the 2006 AVI workshop on BEyond time and errors: novel evaluation methods for information visualization*, pages 1–7. ACM, 2006. (Cited on pages 79, 84, and 90.)
- [123] C. Faloutsos and V. Megalooikonomou. On data mining, compression, and kolmogorov complexity. *Data Mining and Knowledge Discovery*, 15(1):3–20, 2007. (Cited on page 47.)
- [124] C Faloutsos, M Ranganathan, and Y Manolopoulos. Fast subsequence matching in time-series databases. *SIGMOD Record*, 23:419 – 429, 1994. (Cited on pages 22, 36, 43, and 45.)
- [125] L.A. Farwell and E. Donchin. Talking off the top of your head: toward a mental prosthesis utilizing event-related brain potentials. *Electroencephalography and clinical Neurophysiology*, 70(6):510–523, 1988. (Cited on page 218.)
- [126] S.Z. Fatemian, F. Agrafioti, and D. Hatzinakos. Heartid: Cardiac biometric recognition. In *Biometrics: Theory Applications and Systems (BTAS), 2010 Fourth IEEE International Conference on*, pages 1–5. IEEE, 2010. (Cited on page 238.)
- [127] P. Ferreira, P. Azevedo, C. Silva, and R. Brito. Mining approximate motifs in time series. In *Lecture Notes in Computer Science*, volume 4265, pages 89–101. Springer, 2006. (Cited on page 33.)
- [128] J. Figueira, S. Greco, and M. Ehrgott. *Multiple criteria decision analysis: state of the art surveys*. Springer Verlag, 2005. ISBN 038723067X. (Cited on page 49.)
- [129] J.A. Flanagan. A non-parametric approach to unsupervised learning and clustering of symbol strings and sequences. In *Proceedings of the 4th Workshop on Self-Organizing Maps (WSOM03)*, pages 128–133, 2003. (Cited on page 42.)
- [130] C.M. Fonseca, L. Paquete, and M. López-Ibáñez. An improved dimension-sweep algorithm for the hypervolume indicator. In *IEEE Congress on Evolutionary Computation (CEC'2006)*, pages 1157–1163, 2006. (Cited on page 121.)
- [131] A. Frank and A. Asuncion. UCI machine learning repository, 2010. URL <http://archive.ics.uci.edu/ml>. (Cited on pages 130, 216, 217, 222, 223, 226, 227, 229, 230, 231, 232, 234, 235, and 236.)
- [132] A.L. Freire, G.A. Barreto, M. Veloso, and A.T. Varela. Short-term memory mechanisms in neural network learning of robot navigation tasks: a case study. In *Proceedings of the 6th Latin American Robotics Symposium (LARS2009)*, 2009. (Cited on pages 130 and 236.)
- [133] S. Frenkel-Toledo, N. Giladi, C. Peretz, T. Herman, L. Gruendlinger, and J.M. Hausdorff. Treadmill walking as an external pacemaker to improve gait rhythm and stability in parkinson’s disease. *Movement disorders*, 20(9):1109–1114, 2005. (Cited on page 225.)

- [134] E. Frentzos, K. Gratsias, and Y. Theodoridis. Index-based most similar trajectory search. In *IEEE 23rd International Conference on Data Engineering, 2007. ICDE 2007*, pages 816–825, 2007. (Cited on page 41.)
- [135] S. Fröhwrth-Schnatter and S. Kaufmann. Model-based clustering of multiple time series. *Journal of Business and Economic Statistics*, 26(1):78–89, 2008. (Cited on page 27.)
- [136] E. Frøkjær, M. Hertzum, and K. Hornbæk. Measuring usability: are effectiveness, efficiency, and satisfaction really correlated? In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 345–352. ACM, 2000. (Cited on pages 79 and 84.)
- [137] A.W. Fu, E. Keogh, L.Y. Lau, C.A. Ratanamahatana, and R.C.W. Wong. Scaling and time warping in time series querying. *The VLDB Journal - The International Journal on Very Large Data Bases*, 17(4):921, 2008. (Cited on page 40.)
- [138] E. Fuchs, T. Gruber, H. Pree, and B. Sick. Temporal data mining using shape space representations of time series. *Neurocomputing*, 74(1-3):379–393, 2010. (Cited on page 37.)
- [139] M.M. Gaber, A. Zaslavsky, and S. Krishnaswamy. Mining data streams: a review. *ACM Sigmod Record*, 34(2):18–26, 2005. (Cited on page 46.)
- [140] S. Gaffney and P. Smyth. Trajectory clustering with mixtures of regression models. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 63–72. ACM, 1999. (Cited on pages 27 and 42.)
- [141] X. Ge and P. Smyth. Deformable Markov model templates for time-series pattern matching. In *Proceedings of the 6th ACM International conference on Knowledge Discovery and Data Mining*, pages 81–90, 2000. (Cited on page 42.)
- [142] P. Geurts. Pattern extraction for time series classification. In *Proceedings of the 5th European conference on principles of data mining and knowledge discovery*, pages 115 – 127, Freiburg, Germany, 2001. (Cited on page 28.)
- [143] A. Ghias, J. Logan, D. Chamberlin, and B.C. Smith. Query by humming: musical information retrieval in an audio database. In *Proceedings of the third ACM international conference on Multimedia*, pages 231–236. ACM, 1995. (Cited on pages 70 and 166.)
- [144] R. Giot, M. El-Abed, and C. Rosenberger. Greyc keystroke: a benchmark for keystroke dynamics biometric systems. In *Biometrics: Theory, Applications, and Systems, 2009. BTAS'09. IEEE 3rd International Conference on*, pages 1–6. IEEE, 2009. (Cited on page 163.)
- [145] L. Golab and M.T. Ozsü. Issues in data stream management. *ACM Sigmod Record*, 32(2):5–14, 2003. (Cited on page 46.)
- [146] A.L. Goldberger, L.A.N. Amaral, L. Glass, J.M. Hausdorff, P.C. Ivanov, R.G. Mark, J.E. Mietus, G.B. Moody, C.K. Peng, and H.E. Stanley. Physiobank, physiokit, and physionet: Components of a new research resource for complex physiologic signals. *Circulation*, 101(23):e215–e220, 2000. (Cited on pages 221, 225, 231, 234, and 235.)

- [147] D. Goldin and P. Kanellakis. On similarity queries for time-series data: Constraint specification and implementation. In *Principles and Practice of Constraint Programming - CP95*, pages 137–153. Springer, 1995. (Cited on pages 26, 34, and 39.)
- [148] D.Q. Goldin, T.D. Millstein, and A. Kutlu. Bounded similarity querying for time-series data. *Information and Computation*, 194(2):203–241, 2004. (Cited on page 34.)
- [149] JA Gómez-Limón and J. Berbel. Multicriteria analysis of derived water demand functions: a Spanish case study. *Agricultural Systems*, 63(1):49–72, 2000. ISSN 0308-521X. (Cited on page 53.)
- [150] J.W. Gordon. The perceptual attack time of musical tones. *The Journal of the Acoustical Society of America*, 82:88, 1987. (Cited on page 9.)
- [151] D.O. Gorodnichy. Evolution and evaluation of biometric systems. In *Computational Intelligence for Security and Defense Applications, 2009. CISDA 2009. IEEE Symposium on*, pages 1–8. IEEE, 2009. (Cited on pages 155 and 156.)
- [152] M Goto, H Hashiguchi, T Nishimura, and R Oka. Rwc music database : music genre database and musical instrument sound database. In *Proceedings of the 4th International Conference on Music Information Retrieval*, pages 229 – 230, Washington, USA, 2003. (Cited on page 66.)
- [153] S.D. Greenwald, P. Albrecht, G.B. Moody, and R.G. Mark. Estimating confidence limits for arrhythmia detector performance. *Computers in Cardiology*, 12:383–386, 1985. (Cited on page 235.)
- [154] J.M. Grey. Multidimensional perceptual scaling of musical timbres. *Journal of the Acoustical Society of America*, 61(5):1270–1277, 1977. (Cited on page 14.)
- [155] S. Gudmundsson, T.P. Runarsson, and S. Sigurdsson. Support vector machines and dynamic time warping for time series. In *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on*, pages 2772–2776. IEEE, 2008. (Cited on pages 120, 126, and 133.)
- [156] S. Gudmundsson, T.P. Runarsson, and S. Sigurdsson. Test-retest reliability and feature selection in physiological time series classification. *Computer Methods and Programs in Biomedicine*, 105(1):50–60, 2012. (Cited on pages 120, 126, 130, and 231.)
- [157] C. Guger, A. Schlogl, C. Neuper, D. Waltersbacher, T. Strein, and G. Pfurtscheller. Rapid prototyping of an eeg-based brain-computer interface (bci). *Neural Systems and Rehabilitation Engineering, IEEE Transactions on*, 9(1):49–58, 2001. (Cited on page 221.)
- [158] C. Guger, W. Domej, G. Lindner, K. Pfurtscheller, G. Pfurtscheller, and G. Edlinger. Effects of a fast cable car ascent to an altitude of 2700 meters on eeg and ecg. *Neuroscience letters*, 377(1):53–58, 2005. (Cited on pages 130 and 222.)
- [159] F. Gullo, G. Ponti, A. Tagarelli, and S. Greco. A time series representation model for accurate and fast similarity detection. *Pattern Recognition*, 42(11):2998–3014, 2009. (Cited on page 37.)

- [160] G. Guo and S.Z. Li. Content-based audio classification and retrieval by support vector machines. *Neural Networks, IEEE Transactions on*, 14(1):209–215, 2003. (Cited on pages 71, 166, 167, and 173.)
- [161] G. Guo, H.J. Zhang, and S.Z. Li. Boosting for content-based audio classification and retrieval: an evaluation. In *IEEE International Conference on Multimedia and Expo (ICME 2001)*, pages 997–1000, 2001. (Cited on pages 166, 167, and 173.)
- [162] S. Gupta, A. Ray, and E. Keller. Symbolic time series analysis of ultrasonic data for early detection of fatigue damage. *Mechanical Systems and Signal Processing*, 21(2):866–884, 2007. (Cited on page 32.)
- [163] D. Gusfield. *Algorithms on strings, trees, and sequences: computer science and computational biology*. Cambridge Univ Pr, 1997. (Cited on page 45.)
- [164] N. Hammami and M. Bedda. Improved tree model for arabic speech recognition. In *Computer Science and Information Technology (ICCSIT), 2010 3rd IEEE International Conference on*, volume 5, pages 521–526. IEEE, 2010. (Cited on page 216.)
- [165] N. Hammami and M. Sellam. Tree distribution classifier for automatic spoken arabic digit recognition. In *Internet Technology and Secured Transactions, 2009. IC-ITST 2009. International Conference for*, pages 1–4. IEEE, 2009. (Cited on pages 130 and 216.)
- [166] J. Han and M. Kamber. *Data mining: concepts and techniques*. Morgan Kaufmann, 2006. ISBN 1558609016. (Cited on page 27.)
- [167] R.I.D. Harris and R. Sollis. *Applied time series modelling and forecasting*. J. Wiley, 2003. ISBN 0470844434. (Cited on page 30.)
- [168] M. Hassenzahl and D. Ullrich. To do or not to do: Differences in user experience and retrospective judgments depending on the presence or absence of instrumental goals. *Interacting with Computers*, 19(4):429–437, 2007. (Cited on pages 84 and 85.)
- [169] K. Hayashi. Multicriteria analysis for agricultural resource management : a critical survey and future perspectives. *European Journal of Operational Research*, 122(2):486–500, 2000. ISSN 0377-2217. (Cited on page 53.)
- [170] G. Hebrail and B. Hugueney. Symbolic representation of long time-series. In *Symbolic Data Analysis at the 4th European Conference on Principles of Data Mining and Knowledge Discovery*, pages 56–65, 2000. (Cited on page 27.)
- [171] M. Helen and T. Lahti. Query by example methods for audio signals. In *Proceedings of the 7th Nordic Signal Processing Symposium*, pages 302–305, 2006. (Cited on pages 167 and 173.)
- [172] M. Helen and T. Lahti. Query by example in large databases using key-sample distance transformation and clustering. In *Proceedings of the 9th IEEE International Symposium on Multimedia Workshops (ISMW'07)*, pages 303–308, 2007. (Cited on page 71.)
- [173] M. Helen and T. Virtanen. Query by example of audio signals using euclidean distance between gaussian mixture models. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, volume 1, pages 225–228, 2007. (Cited on pages 71 and 166.)

- [174] J.M. Hellerstein, E. Koutsoupas, and C.H. Papadimitriou. On the analysis of indexing schemes. In *Proceedings of the 16th ACM Symposium on Principles of Database Systems*, pages 249–256, 1997. (Cited on page 46.)
- [175] H.L.F. Helmholtz and A.J. Ellis. The sensations of tone: As a physiological basis for the theory of music. 1875. (Cited on page 8.)
- [176] LJ Herrera, H. Pomares, I. Rojas, A. Guillén, A. Prieto, and O. Valenzuela. Recursive prediction for long term time series forecasting using advanced models. *Neurocomputing*, 70(16-18):2870–2880, 2007. (Cited on page 31.)
- [177] J. Himberg, K. Korpiaho, J. Tikanmaki, and H.T.T. Toivonen. Time series segmentation for context recognition in mobile devices. In *Proceedings of the 1st IEEE International Conference on Data Mining*, pages 203–210, 2001. (Cited on page 30.)
- [178] J. Himberg, J. Mantyjarvi, and P. Korpipaa. Using PCA and ICA for exploratory data analysis in situation awareness. In *Proceedings of the International Conference on Multisensor Fusion and Integration for Intelligent Systems*, pages 127–131, 2001. (Cited on page 34.)
- [179] U. Hoffmann, J.M. Vesin, T. Ebrahimi, and K. Diserens. An efficient p300-based brain-computer interface for disabled subjects. *Journal of Neuroscience methods*, 167(1):115–125, 2008. (Cited on pages 130 and 224.)
- [180] C. Holland and O.V. Komogortsev. Biometric identification via eye movement scanpaths in reading. In *Biometrics (IJCB), 2011 International Joint Conference on*, pages 1–8. IEEE, 2011. (Cited on page 163.)
- [181] Y.W. Huang and P.S. Yu. Adaptive query processing for time-series data. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 282–286. ACM, 1999. (Cited on page 36.)
- [182] G. Hulten, L. Spencer, and P. Domingos. Mining time-changing data streams. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 97–106. ACM, 2001. (Cited on page 46.)
- [183] T.A. Hummel. Simulation of human voice timbre by orchestration of acoustic music instruments. In *Proceedings of International Computer Music Conference (ICMC)*, 2005. (Cited on page 185.)
- [184] H. Huttunen, J.P. Kauppi, and J. Tohka. Regularized logistic regression for mind reading with parallel validation. *Proceedings of the ICANN'2011 conference*, page 20, 2011. (Cited on page 228.)
- [185] Y. Ichimaru and GB Moody. Development of the polysomnographic database on cd-rom. *Psychiatry and Clinical Neurosciences*, 53(2):175–177, 1999. (Cited on page 234.)
- [186] P. Indyk, N. Koudas, and S. Muthukrishnan. Identifying representative trends in massive time series data sets using sketches. In *Proceedings of the 26th International Conference on Very Large Data Bases*, pages 363–372. Morgan Kaufmann Publishers Inc., 2000. (Cited on pages 25 and 36.)

- [187] K.T.U. Islam, K. Hasan, Y.K. Lee, and S. Lee. Enhanced 1-nn time series classification using badness of records. In *Proceedings of the 2nd international conference on Ubiquitous information management and communication*, pages 108–113. ACM, 2008. (Cited on pages 120 and 126.)
- [188] S.A. Israel, J.M. Irvine, A. Cheng, M.D. Wiederhold, and B.K. Wiederhold. Ecg to identify individuals. *Pattern Recognition*, 38(1):133–142, 2005. (Cited on page 147.)
- [189] P. Iverson and C.L. Krumhansl. Isolating the dynamic attributes of musical timbre. *The Journal of the Acoustical Society of America*, 94:2595, 1993. (Cited on page 9.)
- [190] A. Jain, L. Hong, and S. Pankanti. Biometric identification. *Communications of the ACM*, 43(2):90–98, 2000. (Cited on pages 146 and 164.)
- [191] A.K. Jain, A. Ross, and S. Prabhakar. An introduction to biometric recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 14(1):4–20, 2004. (Cited on page 147.)
- [192] GJ Janacek, AJ Bagnall, and M. Powell. A likelihood ratio distance measure for the similarity between the Fourier transform of time series. *Lecture Notes in Computer Science*, 3518:737–743, 2005. (Cited on page 42.)
- [193] J.S.R. Jang and H.R. Lee. Hierarchical filtering method for content-based music retrieval via acoustic input. In *Proceedings of the ninth ACM international conference on Multimedia*, page 410. ACM, 2001. (Cited on page 166.)
- [194] J.S.R. Jang, C.L. Hsu, and H.R. Lee. Continuous HMM and Its Enhancement for Singing/Humming Query Retrieval. ISMIR 2005, 6th International Conference on Music Information Retrieval, 2005. (Cited on page 166.)
- [195] J. Jasper and K.R. Othman. Feature extraction for human identification based on envelopegram signal analysis of cardiac sounds in time-frequency domain. In *Electronics and Information Engineering (ICEIE), 2010 International Conference On*, volume 2, pages V2–228. IEEE, 2010. (Cited on pages 147 and 238.)
- [196] C. Jeffery. Synthetic lightning emp data, 2005. URL <http://nis-www.lanl.gov/~eads/datasets/emp>. (Cited on pages 130 and 234.)
- [197] S.L. Jeng and Y.T. Huang. Time Series Classification Based on Spectral Analysis. *Communications in Statistics-Simulation and Computation*, 37(1):132–142, 2008. (Cited on page 29.)
- [198] DL Jones. Fathom: a matlab toolbox for multivariate ecological and oceanographic data analysis, 2002. (Cited on page 100.)
- [199] M.W. Kadous. Grasp: Recognition of australian sign language using instrumented gloves, 1995. (Cited on pages 130 and 217.)
- [200] M.W. Kadous. *Temporal Classification: Extending the Classification Paradigm to Multivariate Time Series*. PhD thesis, The University of New South Wales, 2002. (Cited on pages 130 and 217.)
- [201] M. Käki and A. Aula. Controlling the complexity in comparing search user interfaces via user studies. *Information processing & management*, 44(1):82–91, 2008. (Cited on pages 84, 86, and 87.)

- [202] K. Kalpakis, D. Gada, and V. Puttagunta. Distance measures for effective clustering of ARIMA time-series. In *Proceedings of the IEEE International Conference on Data Mining*, pages 273–280, 2001. (Cited on page 37.)
- [203] B. Kaluza, V. Mirchevska, E. Dovgan, M. Lustrek, and M. Gams. An agent-based approach to care in independent living. *Ambient Intelligence*, pages 177–186, 2010. (Cited on pages 130 and 230.)
- [204] S. Kanade, D. Petrovska-Delacrétaz, and B. Dorizzi. Cancelable iris biometrics and using error correcting codes to reduce variability in biometric data. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 120–127. IEEE, 2009. (Cited on page 163.)
- [205] H. Kauppinen, T. Seppanen, and M. Pietikainen. An experimental comparison of autoregressive and Fourier-based descriptors in 2D shape classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(2):201–207, 1995. (Cited on page 48.)
- [206] J.P. Keener. The perron-frobenius theorem and the ranking of football teams. *SIAM review*, pages 80–93, 1993. (Cited on page 95.)
- [207] A. Kehagias. A hidden Markov model segmentation procedure for hydrological and environmental time series. *Stochastic Environmental Research and Risk Assessment*, 18(2):117–130, 2004. (Cited on page 30.)
- [208] E Keogh and S Kasetty. On the need for time series data mining benchmarks : a survey and empirical demonstration. In *Proceedings of the 8th ACM SIGKDD International conference on knowledge discovery and data mining*, pages 102 – 111, Edmonton, Alberta, Canada, 2002. (Cited on pages 48 and 129.)
- [209] E. Keogh and S. Kasetty. On the need for time series data mining benchmarks: A survey and empirical demonstration. *Data Mining and Knowledge Discovery*, 7(4): 349–371, 2003. (Cited on pages 34, 38, 39, 43, 131, and 132.)
- [210] E. Keogh and M. Pazzani. An enhanced representation of time series which allows fast and accurate classification, clustering and relevance feedback. In *Proceedings of the 4th International Conference of Knowledge Discovery and Data Mining*, pages 239–241. AAAI Press, 1998. (Cited on pages 28, 30, 36, 39, and 47.)
- [211] E. Keogh and C.A. Ratanamahatana. Exact indexing of dynamic time warping. *Knowledge and Information Systems*, 7(3):358–386, 2005. (Cited on pages 26, 40, and 46.)
- [212] E Keogh, K Chakrabarti, and M Pazzani. Locally adaptive dimensionality reduction for indexing large time series databases. In *Proceedings of ACM conference on management of data*, pages 151 – 162, 2001. (Cited on pages 26, 34, and 36.)
- [213] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and Information Systems*, 3(3):263–286, 2001. (Cited on pages 25, 36, and 45.)
- [214] E. Keogh, S. Chu, D. Hart, and M. Pazzani. Segmenting time series: A survey and novel approach. *Data mining in time series databases*, pages 1–21, 2003. (Cited on pages 22, 29, and 46.)

- [215] E. Keogh, J. Lin, and W. Truppel. Clustering of time series subsequences is meaningless: implications for previous and future research. In *3rd IEEE International Conference on Data Mining*, pages 115–122, 2003. (Cited on pages 28 and 33.)
- [216] E. Keogh, S. Lonardi, and C.A. Ratanamahatana. Towards parameter-free data mining. In *Proceedings of 10th ACM international conference on Knowledge discovery and data mining*, pages 206–215, 2004. (Cited on pages 25, 34, 36, 43, and 47.)
- [217] E. Keogh, J. Lin, S.H. Lee, and H.V. Herle. Finding the most unusual time series subsequence: algorithms and applications. *Knowledge and Information Systems*, 11(1):1–27, 2007. (Cited on page 32.)
- [218] G. Kerr, H.J. Ruskin, M. Crane, and P. Doolan. Techniques for clustering gene expression data. *Computers in Biology and Medicine*, 38(3):283–293, 2008. (Cited on page 28.)
- [219] G.A. Kiker, T.S. Bridges, A. Varghese, T.P. Seager, and I. Linkov. Application of multicriteria decision analysis in environmental decision making. *Integrated Environmental Assessment and Management*, 1(2):95–108, 2005. ISSN 1551-3793. (Cited on page 53.)
- [220] S.W. Kim, S. Park, and W.W. Chu. An Index-Based Approach for Similarity Search Supporting Time Warping in Large Sequence Databases. In *Proceedings of the 17th International Conference on Data Engineering*, pages 607–614. IEEE Computer Society, 2001. (Cited on page 26.)
- [221] A. Klami, P. Ramkumar, S. Virtanen, L. Parkkonen, R. Hari, and S. Kaski. Ican-n/pascal2 challenge: Meg mind-reading—overview and results. *Proceedings of the ICANN 2011 Conference*, page 3, 2011. (Cited on pages 130 and 228.)
- [222] C. Kœchlin. *Traité de l’orchestration en quatre volumes*, volume 1. Éditions Max Eschig, 1954. (Cited on pages 2 and 3.)
- [223] E. Koliopoulou. Etude de la description morphologique des sons environnementaux. Technical report, IRCAM, Paris, 2012. (Cited on pages 174, 177, and 178.)
- [224] M. Kontaki, A.N. Papadopoulos, and Y. Manolopoulos. Adaptive similarity search in streaming time series with sliding windows. *Data & Knowledge Engineering*, 63(2):478–502, 2007. (Cited on pages 26 and 46.)
- [225] M. Kontaki, A. Papadopoulos, and Y. Manolopoulos. Similarity Search in Time Series. *Handbook of Research on Innovations in Database Technologies and Applications*, pages 288–299, 2009. (Cited on page 26.)
- [226] F. Korn, HV Jagadish, and C. Faloutsos. Efficiently supporting ad hoc queries in large datasets of time sequences. In *Proceedings of the 1997 ACM SIGMOD international conference on Management of data*, pages 289–300. ACM, 1997. (Cited on pages 25 and 36.)
- [227] T. Koskela. *Neural network methods in analysing and modelling time varying processes*. PhD thesis, Helsinki University of Technology Laboratory of Computational Engineering, 2003. (Cited on page 31.)
- [228] V. Krasteva and I. Jekova. Assessment of ecg frequency and morphology parameters for automatic classification of life-threatening cardiac arrhythmias. *Physiological measurement*, 26:707, 2005. (Cited on pages 130 and 235.)

- [229] M. Kudo, J. Toyama, and M. Shimbo. Multidimensional curve classification using passing-through regions. *Pattern Recognition Letters*, 20(11-13):1103–1111, 1999. (Cited on pages 130 and 227.)
- [230] N. Kumar, N. Lolla, E. Keogh, S. Lonardi, C. Ratanamahatana, and L. Wei. Time-series bitmaps: a practical visualization tool for working with large time series databases. In *SIAM 2005 Data Mining Conference*, pages 531–535, 2005. (Cited on page 37.)
- [231] R.T.H. Laennec. *De l'auscultation médiate ou Traité du diagnostic des maladies des poumons et du coeur, fondé principalement sur ce nouveau moyen d'exploration*, volume 2. Brosson et Chaude, 1819. (Cited on page 236.)
- [232] T.N. Lal, T. Hinterberger, G. Widman, M. Schroder, J. Hill, W. Rosenstiel, C. Elger, B. Scholkopf, and N. Birbaumer. Methods towards invasive human brain computer interfaces. *Advances in Neural Information Processing Systems (NIPS)*, 2005. (Cited on page 218.)
- [233] L.J. Latecki, V. Megalooikonomou, Q. Wang, R. Lakaemper, CA Ratanamahatana, and E. Keogh. Elastic partial matching of time series. *Knowledge Discovery in Databases*, pages 577–584, 2005. (Cited on page 41.)
- [234] L.J. Latecki, Q. Wang, S. Koknar-Tezel, and V. Megalooikonomou. Optimal subsequence bijection. In *IEEE Int. Conf. on Data Mining (ICDM)*, pages 565–570, Omaha, USA, 2007. (Cited on page 41.)
- [235] MH Law and JT Kwok. Rival penalized competitive learning for model-based sequence clustering. In *Pattern Recognition, 2000. Proceedings. 15th International Conference on*, volume 2, pages 2186–2195, 2000. (Cited on page 42.)
- [236] S.H. Lee and J.S. Lim. Parkinson's disease classification using gait characteristics and wavelet-based feature extraction. *Expert Systems with Applications: An International Journal*, 39(8):7338–7344, 2012. (Cited on pages 130 and 225.)
- [237] S. Lemm, C. Schafer, and G. Curio. Bci competition 2003-data set iii: probabilistic modeling of sensorimotor mu rhythms for classification of imaginary hand movements. *Biomedical Engineering, IEEE Transactions on*, 51(6):1077–1080, 2004. (Cited on pages 130 and 219.)
- [238] M. Lesaffre, M. Leman, K. Tanghe, B. De Baets, H. De Meyer, and J.P. Martens. User-dependent taxonomy of musical features as a conceptual framework for musical audio-mining technology. In *Proc. of the Stockholm Music Acoustics Conference*, 2003. (Cited on pages 72 and 84.)
- [239] C.S. Li, P.S. Yu, and V. Castelli. MALM: a framework for mining sequence database at multiple abstraction levels. In *Proceedings of the seventh international conference on Information and knowledge management*, pages 267–272. ACM, 1998. (Cited on page 29.)
- [240] G. Li and A.A. Khokhar. Content-based indexing and retrieval of audio data using wavelets. In *IEEE International Conference on Multimedia and Expo*, volume 2, pages 885–888. IEEE, 2000. (Cited on pages 71 and 166.)

- [241] S.Z. Li. Content-based audio classification and retrieval using the nearest feature line method. *Speech and Audio Processing, IEEE Transactions on*, 8(5):619–625, 2000. (Cited on pages 71, 166, 167, and 173.)
- [242] Y. Li and Y. Hou. Search audio data with the wavelet pyramidal algorithm. *Information processing letters*, 91(1):49–55, 2004. (Cited on page 166.)
- [243] X. Lian and L. Chen. Efficient similarity search over future stream time series. *IEEE Transactions on Knowledge and Data Engineering*, 20(1):40–54, 2007. (Cited on pages 26 and 46.)
- [244] X. Lian, L. Chen, and B. Wang. Approximate similarity search over multiple stream time series. *Lecture Notes in Computer Science*, 4443:962–968, 2010. (Cited on page 26.)
- [245] T.W. Liao. Clustering of time series data—a survey. *Pattern Recognition*, 38(11):1857–1874, 2005. (Cited on page 27.)
- [246] AWC Liew, SH Leung, and WH Lau. Fuzzy image clustering incorporating spatial continuity. *IEEE Proceedings on Vision, Image and Signal Processing*, 147(2):185–192, 2000. (Cited on page 48.)
- [247] J. Lin and E. Keogh. Clustering of time-series subsequences is meaningless: implications for previous and future research. *Knowledge and information systems*, 8(2):154–177, 2005. (Cited on pages 22, 34, and 36.)
- [248] J. Lin and Y. Li. Finding structural similarity in time series data using bag-of-patterns representation. In *Scientific and Statistical Database Management: 21st International Conference, SSDBM 2009 New Orleans, La, USA, June 2-4, 2009 Proceedings*, pages 461–477. Springer, 2009. (Cited on page 42.)
- [249] J. Lin, E. Keogh, S. Lonardi, and B. Chiu. A symbolic representation of time series, with implications for streaming algorithms. In *Proceedings of the 8th ACM SIGMOD workshop on Research issues in data mining and knowledge discovery*, pages 2–11. ACM New York, NY, USA, 2003. (Cited on pages 26, 37, and 73.)
- [250] J. Lin, E. Keogh, S. Lonardi, J.P. Lankford, and D.M. Nystrom. Visually mining and monitoring massive time series. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 460–469. ACM, 2004. (Cited on page 22.)
- [251] J. Lin, E. Keogh, L. Wei, and S. Lonardi. Experiencing sax: a novel symbolic representation of time series. *Data mining and knowledge discovery*, 15(2):107–144, 2007. (Cited on page 61.)
- [252] T. Lin, N. Kaminski, and Z. Bar-Joseph. Alignment and classification of time series gene expression in clinical studies. *Bioinformatics*, 24(13):147–155, 2008. (Cited on pages 22 and 29.)
- [253] Z. Liu, J.X. Yu, X. Lin, H. Lu, and W. Wang. *Locating motifs in time-series data*, pages 343–353. Springer, 2005. (Cited on page 33.)
- [254] F. Lotte, M. Congedo, A. Lécuyer, F. Lamarche, and B. Arnaldi. A review of classification algorithms for EEG-based brain–computer interfaces. *Journal of Neural Engineering*, 4:1–13, 2007. (Cited on page 29.)

- [255] T. Lowitz, M. Ebert, W. Meyer, and B. Hensel. Hidden Markov Models for Classification of Heart Rate Variability in RR Time Series. In *World Congress on Medical Physics and Biomedical Engineering*, pages 1980–1983, Munich, Germany, 2009. Springer. (Cited on page 29.)
- [256] A.A. Luisada, D.M. MacCanon, S. Kumar, and L.P. Feigen. Changing views on the mechanism of the first and second heart sounds. *American Heart Journal*, 88(4):503–514, 1974. (Cited on page 147.)
- [257] J. Ma and S. Perkins. Online novelty detection on temporal sequences. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 613–618. ACM, 2003. (Cited on page 32.)
- [258] H. Mannila and JK Seppnen. Recognizing similar situations from event sequences. In *First SIAM Conference on Data Mining*, pages 1–16, Chicago, IL, USA, 2001. (Cited on page 42.)
- [259] A.J. Mansfield and J.L. Wayman. Best practices in testing and reporting performance of biometric devices. *NPL Report CMSC*, 14(02), 2002. (Cited on page 154.)
- [260] S. Mardle and S. Pascoe. A review of applications of multiple-criteria decision-making techniques to fisheries. *Marine Resource Economics*, 14:41–64, 1999. ISSN 0738-1360. (Cited on page 53.)
- [261] P.F. Marteau. Time warp edit distance with stiffness adjustment for time series matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2): 306–318, 2008. (Cited on page 41.)
- [262] D. Mazzoni and R.B. Dannenberg. Melody matching directly from audio. In *2nd Annual International Symposium on Music Information Retrieval*, pages 17–18, 2001. (Cited on page 166.)
- [263] S. McAdams and A. Bregman. Hearing musical streams. *Computer Music Journal*, 3(4):26–60, 1979. (Cited on page 14.)
- [264] S. McAdams, S. Winsberg, S. Donnadieu, G. Soete, and J. Krimphoff. Perceptual scaling of synthesized musical timbres: Common dimensions, specificities, and latent subject classes. *Psychological research*, 58(3):177–192, 1995. (Cited on pages 14 and 71.)
- [265] V. Megalooikonomou, G. Li, and Q. Wang. A dimensionality reduction technique for efficient similarity analysis of time series databases. In *Proceedings of the thirteenth ACM international conference on Information and knowledge management*, pages 160–161, Washington, D.C., USA, 2004. ACM. (Cited on page 36.)
- [266] V. Megalooikonomou, Q. Wang, G. Li, and C. Faloutsos. A multiresolution symbolic representation of time series. In *Proceedings. 21st International Conference on Data Engineering*, pages 668–679, 2005. (Cited on pages 36 and 42.)
- [267] GA Mendoza and H. Martins. Multi-criteria decision analysis in natural resource management: a critical review of methods and new modelling paradigms. *Forest Ecology and Management*, 230(1-3):1–22, 2006. ISSN 0378-1127. (Cited on pages 49 and 53.)

- [268] I. Mierswa and M. Wurst. Efficient case based feature construction. *Machine Learning: ECML 2005*, pages 641–648, 2005. (Cited on page [170](#).)
- [269] J.R. Miller and E.C. Carterette. Perceptual space for musical structures. *The Journal of the Acoustical Society of America*, 58:711, 1975. (Cited on page [14](#).)
- [270] PG Mills, RF Chamusco, S. Moos, and E. Craige. Echophonocardiographic studies of the contribution of the atrioventricular valves to the first heart sound. *Circulation*, 54(6):944–951, 1976. (Cited on page [147](#).)
- [271] M. Misaki, Y. Kim, P.A. Bandettini, and N. Kriegeskorte. Comparison of multivariate classifiers and response normalizations for pattern-information fmri. *Neuroimage*, 53(1):103–118, 2010. (Cited on pages [129](#) and [131](#).)
- [272] N. Misdariis, A. Minard, P. Susini, G. Lemaitre, S. McAdams, and E. Parizet. Environmental sound perception: Metadescription and modeling based on independent primary studies. *EURASIP Journal on Audio, Speech, and Music Processing*, 2010:26, 2010. (Cited on page [175](#).)
- [273] B. Mjaaland, P. Bours, and D. Gligoroski. Nisk2009-gait mimicking-attack resistance testing of gait authentication systems. *Norsk informasjonssikkerhetskonferanse (NISK)*, 2009. (Cited on page [163](#).)
- [274] Y. Mohammad and T. Nishida. Constrained Motif Discovery in Time Series. *New Generation Computing*, 27(4):319–346, 2009. (Cited on page [33](#).)
- [275] F. Morchen, A. Ultsch, M. Thies, and I. Lohken. Modeling timbre distance with temporal statistics from polyphonic music. *Audio, Speech, and Language Processing, IEEE Transactions on*, 14(1):81–90, 2006. (Cited on page [10](#).)
- [276] M.D. Morse and J.M. Patel. An efficient and accurate method for evaluating time series similarity. In *Proceedings of the 2007 ACM international conference on Management of data*, pages 569–580, 2007. (Cited on page [41](#).)
- [277] A. Mueen, E. Keogh, Q. Zhu, S. Cash, and B. Westover. Exact discovery of time series motifs. In *Proceedings of the SIAM International Conference on Data Mining (SDM)*, pages 473–484, 2009. (Cited on page [33](#).)
- [278] M. Muhammad Fuad and P.F. Marteau. Extending the Edit Distance Using Frequencies of Common Characters. In *Proceedings of the 19th International Conference on Database and Expert Systems Applications*, pages 150–157, Turin, Italy, 2008. Springer. (Cited on page [41](#).)
- [279] Y. Nagasaka, K. Shimoda, and N. Fujii. Multidimensional recording (mdr) and data sharing: An ecological open research and educational platform for neuroscience. *PloS one*, 6(7), 2011. (Cited on pages [130](#), [228](#), and [229](#).)
- [280] A. Nanopoulos, R. Alcock, and Y. Manolopoulos. Feature-based classification of time-series data. In *Information processing and technology*, pages 49–61, 2001. (Cited on pages [29](#) and [37](#).)
- [281] Y. Ogras and H. Ferhatosmanoglu. Online summarization of dynamic time series data. *The VLDB Journal - The International Journal on Very Large Data Bases*, 15(1): 84–98, 2006. (Cited on page [30](#).)

- [282] H.H. Otu and K. Sayood. A new sequence distance measure for phylogenetic tree construction. *Bioinformatics*, 19(16):2122–2130, 2003. (Cited on page 43.)
- [283] R. Ouyang, L. Ren, W. Cheng, and C. Zhou. Similarity search and pattern discovery in hydrological time series data mining. *Hydrological Processes*, 24(9): 1198–1210, 2010. (Cited on page 22.)
- [284] T. Palpanas, E. Keogh, V.B. Zordan, D. Gunopulos, and M. Cardle. Indexing large human-motion databases. In *Proceedings of the 13th international conference on Very large data bases*, pages 780–791, 2004. (Cited on page 34.)
- [285] T. Palpanas, M. Vlachos, E. Keogh, D. Gunopulos, and W. Truppel. Online amnesic approximation of streaming time series. In *20th International Conference on data engineering*, pages 338–349, 2004. (Cited on page 37.)
- [286] T. Palpanas, M. Vlachos, E. Keogh, and D. Gunopulos. Streaming time series summarization using user-defined amnesic functions. *IEEE Transactions on Knowledge and Data Engineering*, 20(7):992–1006, 2008. (Cited on page 30.)
- [287] A. Panuccio, M. Bicego, and V. Murino. A Hidden Markov Model-based approach to sequential data clustering. *Lecture Notes in Computer Science*, 2396:734–743, 2002. (Cited on pages 37 and 42.)
- [288] S. Papadimitriou and P. Yu. Optimal multi-scale patterns in time series streams. In *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*, pages 647–658, Chicago, IL, USA, 2006. (Cited on page 36.)
- [289] S. Papadimitriou, J. Sun, and PS Yu. Local correlation tracking in time series. In *Data Mining, 2006. ICDM'06. Sixth International Conference on*, pages 456–465, 2006. (Cited on page 42.)
- [290] V. Pareto. Cours d'Economie Politique, volume I and II. *F. Rouge, Lausanne*, 250, 1896. (Cited on page 49.)
- [291] S. Park, D. Lee, and W.W. Chu. Fast retrieval of similar subsequences in long sequence databases. In *In 3rd IEEE Knowledge and Data Engineering Exchange Workshop*, pages 60–67, 1999. (Cited on page 29.)
- [292] S. Park, W.W. Chu, J. Yoon, and C. Hsu. Efficient searches for similar subsequences of different lengths in sequence databases. In *Proceedings. 16th International Conference on Data Engineering*, pages 23–32, 2000. ISBN 0769505066. (Cited on page 45.)
- [293] P. Patel, E. Keogh, J. Lin, and S. Lonardi. Mining Motifs in Massive Time Series Databases. In *Proceedings of IEEE International Conference on Data Mining (ICDM02)*, pages 370–377, 2002. (Cited on pages 28 and 32.)
- [294] S. Pauws. CubyHum: a fully operational query by humming system. In *Proceedings of ISMIR*, pages 187–196, 2002. (Cited on page 166.)
- [295] G Peeters. A large set of audio features for sound description in the cuidado project. Technical report, IRCAM, Paris, 2004. (Cited on pages xxvi, 9, 73, 74, and 80.)
- [296] G.G. Peeters and E. Deruty. Automatic morphological description of sounds. *Journal of the Acoustical Society of America*, 123(5):3801, 2008. (Cited on page 71.)

- [297] A. Perina, M. Cristani, U. Castellani, and V. Murino. A new generative feature set based on entropy distance for discriminative classification. *Image Analysis and Processing–ICIAP 2009*, pages 199–208, 2009. (Cited on pages 130 and 222.)
- [298] C.S. Perng, H. Wang, S.R. Zhang, and D.S. Parker. Landmarks : a new model for similarity-based pattern querying in time series databases. In *Proceedings of the 16th International Conference on Data Engineering*, pages 33–42, 2000. (Cited on page 36.)
- [299] P. Perny. Multicriteria filtering methods based on concordance and non-discordance principles. *Annals of Operations Research*, 80:137–165, 1998. ISSN 0254-5330. (Cited on page 53.)
- [300] M.H. Pesaran, D. Pettenuzzo, and A. Timmermann. Forecasting time series subject to multiple structural breaks. *Review of Economic Studies*, 73(4):1057–1084, 2006. (Cited on page 31.)
- [301] K. Phua, T.H. Dat, J. Chen, and L. Shue. Human identification using heart sound. In *Second International Workshop on Multimodal User Authentication*, Toulouse, France, 2006. (Cited on pages 147 and 237.)
- [302] K. Phua, J. Chen, T.H. Dat, and L. Shue. Heart sound as a biometric. *Pattern Recognition*, 41(3):906–919, 2008. (Cited on pages 147 and 238.)
- [303] W. Piston. *Orchestration*. WW Norton New York, 1955. (Cited on pages 2 and 5.)
- [304] R. Plomp. Timbre as a multidimensional attribute of complex tones. *Frequency analysis and periodicity detection in hearing*, pages 397–414, 1970. (Cited on page 14.)
- [305] I. Popivanov and R.J. Miller. Similarity search over time-series data using wavelets. In *Proceedings of the International Conference on Data Engineering*, pages 212–224, 2002. (Cited on pages 25 and 36.)
- [306] R.J. Povinelli, M.T. Johnson, A.C. Lindgren, and J. Ye. Time series classification using Gaussian mixture models of reconstructed phase spaces. *IEEE Transactions on Knowledge and Data Engineering*, 16(6):779–783, 2004. (Cited on page 29.)
- [307] F. Prat, A. Marzal, S. Martin, R. Ramos-Garuo, and M.J. Castro. A template-based recognition system for on-line handwritten characters. *Journal of Information Science and Engineering*, 25(3):779–791, 2009. (Cited on pages 130 and 229.)
- [308] J. Pressing. *Synthesizer performance and real-time techniques*. AR Editions, Inc. Madison, WI, USA, 1992. (Cited on page 77.)
- [309] D. Pressnitzer, S. McAdams, S. Winsberg, and J. Fineberg. Perception of musical tension for nontonal orchestral timbres and its relation to psychoacoustic roughness. *Perception and Psychophysics*, 62(1):66–80, 2000. ISSN 0031-5117. (Cited on page 9.)
- [310] H. Proença and L.A. Alexandre. The nice. i: noisy iris challenge evaluation-part i. In *Biometrics: Theory, Applications, and Systems, 2007. BTAS 2007. First IEEE International Conference on*, pages 1–4. IEEE, 2007. (Cited on page 163.)
- [311] D. Psenicka. Sporch: An algorithm for orchestration based on spectral analyses of recorded sounds. In *Proceedings of International Computer Music Conference (ICMC)*, 2003. (Cited on page 184.)

- [312] H. Qi, P. Hartono, K. Suzuki, and S. Hashimoto. Sound database retrieved by sound. *Acoustical Science and Technology*, 23(6):293–300, 2002. (Cited on pages 72 and 166.)
- [313] W. Quigguo, M. Fei, W. Yijun, G. Xiaorong, and G. Shang kai. Feature combination for classifying single-trial ecog during motor imagery of different sessions. *Progress in natural science*, 17(7):851–858, 2007. (Cited on pages 130 and 218.)
- [314] M. Radovanović, A. Nanopoulos, and M. Ivanović. Hubs in space: Popular nearest neighbors in high-dimensional data. *The Journal of Machine Learning Research*, 999:2487–2531, 2010. (Cited on pages 120 and 126.)
- [315] M. Radovanovic, A. Nanopoulos, and M. Ivanovic. Time-series classification in many intrinsic dimensions. In *10th SIAM International Conference on Data Mining*. Citeseer, 2010. (Cited on pages 120 and 126.)
- [316] D. Rafiei and A. Mendelzon. Efficient Retrieval of Similar Time Sequences Using DFT. In *Proceedings. 5th International Conference of Foundations of Data Organization and Algorithms*, pages 249–257, 1998. (Cited on page 25.)
- [317] S. Rahimi, L. Gandy, and N. Mogharreban. A web-based high-performance multicriteria decision support system for medical diagnosis. *International Journal of Intelligent Systems*, 22(10):1083–1099, 2007. ISSN 1098-111X. (Cited on page 53.)
- [318] A. Rakotomamonjy and V. Guigue. Bci competition iii: dataset ii-ensemble of svms for bci p300 speller. *Biomedical Engineering, IEEE Transactions on*, 55(3): 1147–1154, 2008. (Cited on pages 130 and 218.)
- [319] C. Ratanamahatana and D. Wanichsan. Stopping Criterion Selection for Efficient Semi-supervised Time Series Classification. *Studies in Computational Intelligence*, 149:1–14, 2008. (Cited on page 29.)
- [320] C. Ratanamahatana, E. Keogh, A.J. Bagnall, and S. Lonardi. A novel bit level time series representation with implication of similarity search and clustering. *Advances in Knowledge Discovery and Data Mining*, pages 771–777, 2005. (Cited on pages 26 and 37.)
- [321] C.A. Ratanamahatana and E. Keogh. Everything you know about dynamic time warping is wrong. In *Third Workshop on Mining Temporal and Sequential Data*, pages 1–11, Seattle, WA, USA, 2004. (Cited on page 40.)
- [322] C.A. Ratanamahatana and E. Keogh. Making time-series classification more accurate using learned constraints. In *Proceedings of SIAM International Conference on Data Mining*, pages 11–22, 2004. (Cited on page 40.)
- [323] G. Ravi, S.K. Gupta, and MB Ray. Multiobjective optimization of cyclone separators using genetic algorithm. *Ind. Eng. Chem. Res*, 39(11):4272–4286, 2000. (Cited on page 53.)
- [324] KV Ravi Kanth, D. Agrawal, and A. Singh. Dimensionality reduction for similarity searching in dynamic databases. *ACM SIGMOD Record*, 27(2):166–176, 1998. (Cited on page 36.)
- [325] G. Reeves, J. Liu, S. Nath, and F. Zhao. Managing massive time series streams with multi-scale compressed trickles. *Proceedings of the VLDB Endowment*, 2(1): 97–108, 2009. (Cited on pages 34 and 36.)

- [326] T. Rehman and C. Romero. The application of the MCDM paradigm to the management of agricultural systems: some basic considerations. *Agricultural Systems*, 41(3):239–255, 1993. ISSN 0308-521X. (Cited on page 53.)
- [327] G. Reinert, S. Schbath, and M.S. Waterman. Probabilistic and statistical properties of words: an overview. *Journal of Computational Biology*, 7(1-2):1–46, 2000. (Cited on page 42.)
- [328] MJ Reyes-Gomez and DPW Ellis. Selection, parameter estimation, and discriminative training of hidden markov models for general audio modeling. In *Proceedings of the IEEE International Conference on Multimedia and Expo*, volume 1, pages 70–73, 2003. (Cited on pages 167 and 173.)
- [329] J.C. Risset. *Hauteur et timbre des sons*. IRCAM, Paris, 1978. (Cited on pages 8 and 9.)
- [330] JJ Rodriguez and LI Kuncheva. Time series classification: Decision forests and SVM on interval and DTW features. In *Proc Workshop on Time Series Classification, 13th International Conference on Knowledge Discovery and Data mining*, 2007. (Cited on page 29.)
- [331] C. Romero and T. Rehman. Natural resource management and the use of multiple criteria decision-making techniques: A review. *European Review of Agricultural Economics*, 14(1):61, 1987. ISSN 0165-1587. (Cited on page 53.)
- [332] F. Rose and J. Hetrick. Spectral analysis as a resource for contemporary orchestration technique. In *Proceedings of Conference on Interdisciplinary Musicology*, volume 2005, 2005. (Cited on page 184.)
- [333] Y. Sakurai, M. Yoshikawa, S. Uemura, and H. Kojima. The A-tree: An index structure for high-dimensional spaces using relative approximation. In *Proceedings of the 26th International Conference on Very Large Data Bases*, pages 516–526, 2000. (Cited on page 45.)
- [334] Y. Sakurai, M. Yoshikawa, and C. Faloutsos. FTW: fast similarity search under the time warping distance. In *Proceedings of the 24th ACM Symposium on Principles of database systems*, pages 326–337, 2005. (Cited on page 26.)
- [335] S. Salvador and P. Chan. Toward accurate dynamic time warping in linear time and space. *Intelligent Data Analysis*, 11(5):561–580, 2007. (Cited on page 40.)
- [336] S. Salvador, P. Chan, and J. Brodie. Learning states and rules for time series anomaly detection. In *Proc. 17th International FLAIRS Conference*, pages 300–305, 2004. (Cited on page 32.)
- [337] S.L. Salzberg. On comparing classifiers: Pitfalls to avoid and a recommended approach. *Data Mining and knowledge discovery*, 1(3):317–328, 1997. (Cited on pages 129, 131, and 134.)
- [338] S.L. Salzberg. On comparing classifiers: A critique of current research and methods. *Data mining and knowledge discovery*, 1(1), 1999. (Cited on pages 129 and 131.)

- [339] S. Samson, R.J. Zatorre, and J.O. Ramsay. Multidimensional scaling of synthetic musical timbre: Perception of spectral and temporal characteristics. *Canadian Journal of Experimental Psychology*, 51(4):307–315, 1997. ISSN 1196-1961. (Cited on page 14.)
- [340] P. Schaeffer. *Traité des objets musicaux*. Seuil, 1966. (Cited on pages 4, 6, and 174.)
- [341] S. Schliebs, H. Hamed, and N. Kasabov. Reservoir-based evolving spiking neural network for spatio-temporal pattern recognition. In *Neural Information Processing*, pages 160–168. Springer, 2011. (Cited on pages 130 and 227.)
- [342] JF Schouten. The perception of timbre. In *Reports of the 6th International Congress on Acoustics*, volume 76, 1968. (Cited on page 10.)
- [343] T.B. Sebastian, P.N. Klein, and B.B. Kimia. On aligning curves. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(1):116–125, 2003. (Cited on page 48.)
- [344] P. Sebastiani, M. Ramoni, P. Cohen, J. Warwick, and J. Davis. Discovering dynamics using Bayesian clustering. *Lecture Notes in Computer Science*, 1642: 199–209, 1999. (Cited on pages 37 and 42.)
- [345] J. Seibert. Multi-criteria calibration of a conceptual runoff model using a genetic algorithm. *Hydrology and Earth System Sciences*, 4(2):215–224, 2000. (Cited on page 53.)
- [346] J. Seinfeld and W. McBride. Optimization with multiple performance criteria. *Ind. Eng. Chem. Process Des. Develop*, 9(1):53–57, 1970. (Cited on page 53.)
- [347] A. Sfetsos and C. Siriopoulos. Time series forecasting with a hybrid clustering scheme and pattern recognition. *IEEE Transactions on Systems, Man and Cybernetics, Part A*, 34(3):399–405, 2004. (Cited on page 31.)
- [348] X. Shao, C. Xu, and M.S. Kankanhalli. Applying neural network on the content-based audio classification. In *Proceedings of the 4th IEEE Joint International Conference on Information, Communications and Signal Processing*, volume 3, pages 1821–1825, 2003. (Cited on pages 167 and 173.)
- [349] D.E. Shasha and Y. Zhu. *High performance discovery in time series: techniques and case studies*. Springer-Verlag New York Inc, 2004. (Cited on pages 30 and 36.)
- [350] H. Shatkey and SB Zdonik. Approximate queries and representations for large data sequences. In *Data Engineering, 1996. Proceedings of the Twelfth International Conference on*, pages 536–545, 1996. (Cited on pages 25, 29, 36, and 42.)
- [351] J. Shieh and E. Keogh. isax : indexing and mining terabyte sized time series. In *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 623–631. ACM, 2008. (Cited on pages 37, 38, 46, and 62.)
- [352] V. G. Sigillito, S. P. Wing, L. V. Hutton, and K. B. Baker. Classification of radar returns from the ionosphere using neural networks. *Johns Hopkins APL Technical Digest*, pages 262–266, 1989. (Cited on page 226.)
- [353] Y.N. Singh and SK Singh. Evaluation of electrocardiogram for biometric authentication. *Journal of Information Security*, 3:39–48, 2012. (Cited on page 163.)

- [354] M. Slaney. Mixtures of probability experts for audio retrieval and indexing. In *Proceedings of the IEEE International Conference on Multimedia and Expo*, volume 1, pages 345–348. IEEE, 2002. (Cited on pages 71 and 167.)
- [355] M. Slaney. Semantic-audio retrieval. In *Acoustics, Speech, and Signal Processing (ICASSP), 2002 IEEE International Conference on*, volume 4, pages IV–4108, 2002. (Cited on pages 71 and 167.)
- [356] P. Smyth. Clustering sequences with hidden Markov models. *Advances in Neural Information Processing Systems*, pages 648–654, 1997. (Cited on page 27.)
- [357] H. Song and G. Li. Tourism demand modelling and forecasting—A review of recent research. *Tourism Management*, 29(2):203–220, 2008. (Cited on pages 22 and 31.)
- [358] A. Sorjamaa, J. Hao, N. Reyhani, Y. Ji, and A. Lendasse. Methodology for long-term prediction of time series. *Neurocomputing*, 70(16-18):2861–2869, 2007. (Cited on page 31.)
- [359] D. Srisai and C.A. Ratanamahatana. Efficient Time Series Classification under Template Matching Using Time Warping Alignment. In *Proceedings of the Fourth International Conference on Computer Sciences and Convergence Information Technology*, pages 685–690. IEEE, 2009. (Cited on page 29.)
- [360] T. Stiefmeier, D. Roggen, and G. Troster. Gestures are strings: Efficient online gesture spotting and classification using string matching. In *Proceedings of the ICST 2nd international conference on Body area networks*, pages 1–8, Florence, Italy, 2007. (Cited on pages 26, 34, and 36.)
- [361] RG Stockwell, L. Mansinha, and RP Lowe. Localization of the complex spectrum: the s transform. *Signal Processing, IEEE Transactions on*, 44(4):998–1001, 1996. (Cited on pages 151, 238, and 239.)
- [362] ZR Struzik, A. Siebes, and A. CWI. Measuring time series similarity through large singular features revealed with wavelet transformation. In *Proceedings of the Tenth International Workshop on Database and Expert Systems Applications*, pages 162–166, 1999. (Cited on page 36.)
- [363] A. Subasi. EEG signal classification using wavelet feature extraction and a mixture of expert model. *Expert Systems with Applications*, 32(4):1084–1093, 2007. (Cited on page 29.)
- [364] SR Subramanya and A. Youssef. Wavelet-based indexing of audio data in audio/-multimedia databases. In *Multi-Media Database Management Systems, 1998. Proceedings. International Workshop on*, pages 46–53. IEEE, 1998. (Cited on page 166.)
- [365] SR Subramanya, R. Simha, B. Narahari, and A. Youssef. Transform-based indexing of audio data for multimedia databases. In *Proceedings of the IEEE International Conference on Multimedia Computing and Systems' 97*, pages 211–218. IEEE, 1997. (Cited on pages 71 and 166.)
- [366] F. Sufi, I. Khalil, and J. Hu. Ecg-based authentication. *Handbook of Information and Communication Security*, pages 309–331, 2010. (Cited on page 147.)

- [367] S. Sundaram and S. Narayanan. Audio retrieval by latent perceptual indexing. In *Acoustics, Speech and Signal Processing. ICASSP 2008. IEEE International Conference on*, pages 49–52. IEEE, 2008. (Cited on page 71.)
- [368] J. Sundberg. Musical significance of musicians syllable choice in improvised nonsense text singing: A preliminary study. *Phonetica*, 51(1-3):132–145, 1994. (Cited on page 77.)
- [369] J. Sundberg. Level and Center Frequency of the Singer’s Formant* 1. *Journal of voice*, 15(2):176–186, 2001. (Cited on page 77.)
- [370] N. Takama and T. Umeda. Multi-level, multi-objective optimization in process engineering. *Chemical Engineering Science*, 36(1):129–136, 1981. ISSN 0009-2509. (Cited on page 53.)
- [371] P. Tan and D. Dowe. Mml inference of oblique decision trees. *AI 2004: Advances in Artificial Intelligence*, pages 321–338, 2005. (Cited on pages 130 and 234.)
- [372] H. Tang and S.S. Liao. Discovering original motifs with different lengths from time series. *Knowledge-Based Systems*, 21(7):666–671, 2008. (Cited on page 33.)
- [373] D. Tardieu. *Modèles d’instruments pour l’aide à l’orchestration*. PhD thesis, 2008. (Cited on page 185.)
- [374] D. Tardieu and X. Rodet. An instrument timbre model for computer aided orchestration. In *Applications of Signal Processing to Audio and Acoustics, 2007 IEEE Workshop on*, pages 347–350. IEEE, 2007. (Cited on pages 185 and 187.)
- [375] D. Tardieu, G. Carpentier, and X. Rodet. Computer-aided orchestration based on probabilistic instruments models and genetic exploration. In *The Proceedings of International Computer Music Conference*, pages 188–91, 2007. (Cited on page 185.)
- [376] T. Theodoridis and H. Hu. Action classification of 3d human models using dynamic anns for mobile robot surveillance. In *Robotics and Biomimetics, 2007. ROBIO 2007. IEEE International Conference on*, pages 371–376. IEEE, 2007. (Cited on pages 130, 231, and 235.)
- [377] T.A. Traill and N.J. Fortuin. Presystolic mitral closure sound in aortic regurgitation with left ventricular hypertrophy and first degree heart block. *British Heart Journal*, 48(1):78, 1982. (Cited on page 147.)
- [378] D.H. Tran, Y.R. Leng, and H. Li. Feature integration for heart sound biometrics. In *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*, pages 1714–1717. IEEE, 2010. (Cited on page 238.)
- [379] R.S. Tsay. *Analysis of financial time series*. Wiley-Interscience, 2005. ISBN 0471690740. (Cited on page 30.)
- [380] A. Tsoukas, M. Tirrell, and G. Stephanopoulos. Multiobjective dynamic optimization of semibatch copolymerization reactors. *Chemical Engineering Science*, 37(12):1785–1795, 1982. ISSN 0009-2509. (Cited on page 53.)
- [381] A. Uitdenbogerd and J. Zobel. Melodic matching techniques for large music databases. In *Proceedings of the seventh ACM international conference on Multimedia (Part 1)*, page 66. ACM, 1999. (Cited on page 166.)

- [382] M. Van Gerven and O. Jensen. Attention modulations of posterior alpha as a control signal for two-dimensional brain-computer interfaces. *Journal of neuroscience methods*, 179(1):78–84, 2009. (Cited on pages 130 and 220.)
- [383] D.A. van Leeuwen, A.F. Martin, M.A. Przybicki, and J.S. Bouten. Nist and nfi-tno evaluations of automatic speaker recognition. *Computer Speech & Language*, 20(2): 128–158, 2006. (Cited on page 163.)
- [384] A. Van Oosterom, R. Hoekema, and GJH Uijen. Geometrical factors affecting the interindividual variability of the ecg and the vcg. *Journal of electrocardiology*, 33: 219–227, 2000. (Cited on page 164.)
- [385] K.T. Vasko and H.T.T. Toivonen. Estimating the number of segments in time series data using permutation tests. In *Proceedings of the IEEE International Conference on Data Mining*, pages 466–473, 2002. (Cited on page 30.)
- [386] T. Virtanen and M. Helen. Probabilistic model based similarity measures for audio query-by-example. In *Applications of Signal Processing to Audio and Acoustics, 2007 IEEE Workshop on*, pages 82–85, 2007. (Cited on pages 14, 56, 71, and 167.)
- [387] M. Vlachos, D. Gunopoulos, and G. Kollios. Discovering similar multidimensional trajectories. In *Proceedings of the 18th International Conference on Data Engineering*, pages 673–684. IEEE Computer Society, 2002. (Cited on pages 26, 34, and 41.)
- [388] M. Vlachos, J. Lin, E. Keogh, and D. Gunopoulos. A wavelet-based anytime algorithm for k-means clustering of time series. In *Proc. Workshop on Clustering High Dimensionality Data and Its Applications*, pages 23–30, 2003. (Cited on page 27.)
- [389] M. Vlachos, D. Gunopoulos, and G. Das. Indexing time-series under conditions of noise. *Data mining in time series databases*, pages 67–100, 2004. (Cited on pages 26 and 36.)
- [390] M. Vlachos, P. Yu, and V. Castelli. On periodicity detection and structural periodic similarity. In *SIAM International Conference on Data Mining*, pages 449–460, Newport Beach, CA, 2005. (Cited on page 42.)
- [391] M. Vlachos, M. Hadjieleftheriou, D. Gunopoulos, and E. Keogh. Indexing multi-dimensional time-series. *The VLDB Journal*, 15(1):1–20, 2006. (Cited on pages 41 and 46.)
- [392] G. Von Bismarck. Sharpness as an attribute of the timbre of steady sounds. *Acustica*, 30(3):159–172, 1974. (Cited on page 9.)
- [393] N. Wagner, Z. Michalewicz, M. Khouja, and R.R. McGregor. Time series forecasting for dynamic environments: the DyFor genetic program model. *IEEE transactions on evolutionary computation*, 11(4):433–452, 2007. (Cited on page 31.)
- [394] C. Wan and M. Liu. Content-based audio retrieval with relevance feedback. *Pattern recognition letters*, 27(2):85–92, 2006. (Cited on page 166.)
- [395] C. Wan, M. Liu, and L. Wang. Content-based sound retrieval for web application. *Web Intelligence: Research and Development*, pages 389–393, 2001. (Cited on pages 71 and 166.)

- [396] Y. Wang, Z. Zhang, Y. Li, X. Gao, S. Gao, and F. Yang. Bci competition 2003-data set iv: an algorithm based on cssd and fda for classifying single-trial eeg. *Biomedical Engineering, IEEE Transactions on*, 51(6):1081–1086, 2004. (Cited on pages 130 and 219.)
- [397] Y. Wang, F. Agraftioti, D. Hatzinakos, and K.N. Plataniotis. Analysis of human electrocardiogram for biometric recognition. *EURASIP journal on Advances in Signal Processing*, 2008:19, 2008. (Cited on page 147.)
- [398] AS Weigend and NA Gershenfeld. *Time Series Prediction: forecasting the future and understanding the past*. Addison Wesley, 1994. (Cited on page 22.)
- [399] G.M. Weiss. Mining with rarity: a unifying framework. *ACM SIGKDD Explorations Newsletter*, 6(1):7–19, 2004. (Cited on pages 22 and 32.)
- [400] D.L. Wessel. Timbre space as a musical control structure. *Computer music journal*, 3(2):45–52, 1979. ISSN 0148-9267. (Cited on page 14.)
- [401] G. Wichern, J. Xue, H. Thornburg, and A. Spanias. Distortion-aware query-by-example for environmental sounds. In *Applications of Signal Processing to Audio and Acoustics, 2007 IEEE Workshop on*, pages 335–338, 2007. (Cited on page 72.)
- [402] B. Williams, M. Toussaint, and A. Storkey. Extracting motion primitives from natural handwriting data. *Artificial Neural Networks–ICANN 2006*, pages 634–643, 2006. (Cited on page 222.)
- [403] B. Williams, M. Toussaint, and A. Storkey. Modelling motion primitives and their timing in biologically executed movements. *Advances in Neural Information Processing Systems*, 20:1609–1616, 2008. (Cited on page 222.)
- [404] E. Wold, T. Blum, D. Keislar, and J. Wheaton. Content-based classification, search, and retrieval of audio. *Multimedia, IEEE*, 3(3):27–36, 1996. (Cited on pages 71, 166, 167, and 173.)
- [405] X. Xi, E. Keogh, C. Shelton, L. Wei, and C.A. Ratanamahatana. Fast time series classification using numerosity reduction. In *Proceedings of the 23rd international conference on Machine learning ICML 06*, volume 150, pages 1033–1040. ACM Press, 2006. (Cited on pages 29, 126, and 132.)
- [406] X. Xi, E. Keogh, L. Wei, and A. Mafra-neto. Finding Motifs in Database of Shapes. In *Proc. of SIAM International Conference on Data Mining*, pages 249–260, Minneapolis, Minnesota, USA, 2007. (Cited on page 48.)
- [407] H. Xia, G.A. Garcia, J.C. McBride, A. Sullivan, T. De Bock, J. Bains, D.C. Wortham, and X. Zhao. Computer algorithms for evaluating the quality of ecgs in real time. 2012. (Cited on pages 130 and 221.)
- [408] J. Xie and WY Yan. Pattern-based characterization of time series. *International Journal of Information and Systems Science*, 3(3):479–491, 2007. (Cited on page 37.)
- [409] Y. Xiong and D.Y. Yeung. Time series clustering with ARMA mixtures. *Pattern Recognition*, 37(8):1675–1689, 2004. (Cited on page 42.)

- [410] J. Xue, G. Wichern, H. Thornburg, and A. Spanias. Fast query by example of environmental sounds via robust and efficient cluster-based indexing. In *Acoustics, Speech and Signal Processing. ICASSP 2008. IEEE International Conference on*, pages 5–8, 2008. (Cited on page 72.)
- [411] RN Yadav, PK Kalra, and J. John. Time series prediction with single multiplicative neuron model. *Applied soft computing*, 7(4):1157–1163, 2007. (Cited on page 31.)
- [412] N. Yager and T. Dunstone. The biometric menagerie. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 32(2):220–230, 2010. (Cited on page 161.)
- [413] D. Yankov, E. Keogh, J. Medina, B. Chiu, and V. Zordan. Detecting time series motifs under uniform scaling. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 844–853. ACM, 2007. (Cited on page 33.)
- [414] D. Yankov, E. Keogh, and U. Rebbapragada. Disk aware discord discovery: Finding unusual time series in terabyte sized datasets. *Knowledge and Information Systems*, 17(2):241–262, 2008. (Cited on page 32.)
- [415] P.O. Yapo, H.V. Gupta, and S. Sorooshian. Multi-objective global optimization for hydrologic models. *Journal of hydrology*, 204(1-4):83–97, 1998. ISSN 0022-1694. (Cited on page 53.)
- [416] D. Ye, X. Wang, E. Keogh, and A. Mafra-Neto. Autocannibalistic and Anyspace Indexing Algorithms with Applications to Sensor Data Mining. In *The SIAM International Conference on Data Mining (SDM 2009)*, pages 85–96, Sparks, Nevada, 2009. (Cited on page 47.)
- [417] L. Ye and E. Keogh. Time series shapelets: a new primitive for data mining. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 947–956. ACM, 2009. (Cited on pages 37 and 48.)
- [418] B.K. Yi and C. Faloutsos. Fast time sequence indexing for arbitrary Lp norms. In *Proceedings of the 26th International Conference on Very Large Data Bases*, pages 385–394, 2000. (Cited on pages 25, 36, 38, 39, and 45.)
- [419] B.K. Yi, HV Jagadish, and C. Faloutsos. Efficient retrieval of similar time sequences under time warping. In *Data Engineering, 1998. Proceedings., 14th International Conference on*, pages 201–208, 1998. (Cited on page 26.)
- [420] H. Yoon, K. Yang, and C. Shahabi. Feature subset selection and feature ranking for multivariate time series. *IEEE transactions on knowledge and data engineering*, pages 1186–1198, 2005. (Cited on page 27.)
- [421] A. Yoshitaka and T. Ichikawa. A survey on content-based retrieval for multimedia databases. *IEEE Transactions on Knowledge and Data Engineering*, 11(1):81–93, 1999. (Cited on page 70.)
- [422] A. Ypma and R.P.W. Duin. Novelty detection using self-organizing maps. *Progress in Connectionist-Based Information Systems*, 2:1322–1325, 1997. (Cited on page 32.)
- [423] CC Yuen, SK Gupta, and AK Ray. Multi-objective optimization of membrane separation modules using genetic algorithm. *Journal of Membrane Science*, 176(2): 177–196, 2000. ISSN 0376-7388. (Cited on page 53.)

- [424] MF Zafar, D. Mohamad, and MM Anwar. Recognition of online isolated hand-written characters by backpropagation neural nets using sub-character primitive features. In *Multitopic Conference, 2006. INMIC'06. IEEE*, pages 157–162. IEEE, 2006. (Cited on page 222.)
- [425] Y.Y. Zhan, X.Y. Chen, and R.C. Xu. Outlier detection algorithm based on pattern representation of time series. *Application Research of Computers*, 24(11):96–99, 2007. (Cited on page 37.)
- [426] T. Zhang and C.C.J. Kuo. Classification and retrieval of sound effects in audiovisual data management. In *Signals, Systems, and Computers, 1999. Conference Record of the Thirty-Third Asilomar Conference on*, volume 1, pages 730–734. IEEE, 1999. (Cited on pages 71 and 167.)
- [427] T. Zhang and C.C.J. Kuo. Hierarchical classification of audio data for archiving and retrieving. In *Proceedings of the IEEE International Conference on Acoustic, Speech and Signal Processing*, pages 3001–3004, 1999. (Cited on page 71.)
- [428] X. Zhang, J. Wu, X. Yang, H. Ou, and T. Lv. A novel pattern extraction method for time series classification. *Optimization and Engineering*, 10(2):253–271, 2009. (Cited on pages 29 and 33.)
- [429] X.L. Zhang, H. Begleiter, B. Porjesz, W. Wang, and A. Litke. Event related potentials during object recognition tasks. *Brain Research Bulletin*, 38(6):531–538, 1995. (Cited on page 223.)
- [430] Z. Zhang, Y. Shi, and G. Gao. A rough set-based multiple criteria linear programming approach for the medical diagnosis and prognosis. *Expert Systems with Applications*, 36(5):8932–8937, 2009. ISSN 0957-4174. (Cited on page 53.)
- [431] S. Zhong and J. Ghosh. HMMs and coupled HMMs for multi-channel EEG classification. In *Proceedings of the IEEE International Joint Conference on Neural Networks*, pages 1154–1159, 2002. (Cited on pages 29, 130, and 223.)
- [432] S. Zhong, T.M. Khoshgoftaar, and N. Seliya. Clustering-based network intrusion detection. *International Journal of Reliability Quality and Safety Engineering*, 14(2): 169–187, 2007. (Cited on pages 22, 28, and 31.)
- [433] Y. Zhu and D. Shasha. Query by humming: a time series database approach. In *Proc. of SIGMOD*, page 675, 2003. (Cited on page 70.)
- [434] Y. Zhu and D. Shasha. Warping indexes with envelope transforms for query by humming. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*, pages 181–192. ACM, 2003. (Cited on page 166.)
- [435] Y. Zhu, M. Kankanhalli, and C. Xu. Pitch tracking and melody slope matching for song retrieval. *Advances in Multimedia Information Processing (PCM 2001)*, pages 530–537, 2001. (Cited on page 166.)
- [436] E. Zitzler, L. Thiele, M. Laumanns, C.M. Fonseca, and V.G. da Fonseca. Performance assessment of multiobjective optimizers: An analysis and review. *Evolutionary Computation, IEEE Transactions on*, 7(2):117–132, 2003. (Cited on page 121.)