



# *MIRtoolbox 1.8.1*

## *User's Manual*

Olivier Lartillot

Currently: MIRAGE project (funded by Research Council of Norway)

University of Oslo, Norway

RITMO Centre for Interdisciplinary Studies in Rhythm, Time and Motion

[olartillot@gmail.com](mailto:olartillot@gmail.com)

July, 30th, 2021

# TABLE OF CONTENTS

<i>Table of contents</i>	2
<i>1. Introduction</i>	5
Conditions of Use	5
Please Register	5
Documentation and Support	5
Background	6
MIRtoolbox Objectives	7
MIRtoolbox Features	8
Installation	10
Help and demos	12
MIRtoolbox Interface	13
<i>2. Basic Operators</i>	19
miraudio	19
mirframe vs. 'Frame'	23
mirfilterbank	26
mirenvelope	30
mirspectrum	37
mircepstrum	46
mira autocor	50
*	57
mirflux	58
mirsum	63
mirpeaks	65
mirsegment	71
mirplay	74
mirplayer	76
mirsave	77
mirlength	79
mirfdw	80
<i>3. Feature Extractors</i>	82
3.1. Dynamics	82
mirrms	82
mirsegment(..., 'RMS')	85
mirlowenergy	86



3.2. Rhythm	88
mirfluctuation	88
mirbeatspectrum	91
mirevents (previously mironsets)	93
mireventdensity	102
mirtempo	103
mirmetre	109
mirmetroid	112
mirpulseclarity	114
3.3. Timbre	117
mirattacktime	117
mirattackslope	119
mirattackleap	121
mirdecaytime	123
mirdecayleap	125
mirdecayslope (previously mirdecreaseslope)	127
mirduration	129
mirzerocross	130
mirrolloff	132
mirbrightness	134
mircentroid, mirspread, mirskewness, mirkurtosis, mirflatness, mirentropy	136
mirmfcc	137
mirroughness	140
mirregularity	143
3.4. Pitch	145
mirpitch	145
mirmidi	150
mirinharmonic	151
3.5. Tonality	153
mirchromagram	153
mirkeystrength	157
mirkey	160
mirmode	162
mirkeysom	164
mirtonalcentroid	166
mirhcdf	168
mirsegment(..., 'HCDF')	169

## 4. High-level features

170

4.1. Structure and form	170
mirsimatrix	170
mirnovelty	176
mirsegment(..., 'Novelty')	180
4.2. Statistics	182
mirmean	182
mirstd	182
mirmedian	182
mirstat	183
mirhisto	184
mirzerocross	185
mircentroid	186
mirspread	188
mirskewness	189
mirkurtosis	191
mirflatness	192
mirentropy	193
mirfeatures	195
mirmap	198
4.3. Predictions	201
miremotion	201
mirclassify	208
mircluster	209
4.4. Similarity and Retrieval	211
mirdist	211
mirquery	212
4.5. Exportation	213
mirgetdata	213
mirexport	215

## 5. *Advanced use of MIRtoolbox* 217

5.1. Parallel processing	217
5.2. Progressive storage of stats	218
5.3. Interface preferences	219
5.4. get220	
5.5. Memory management	221

## *References* 226

# I. INTRODUCTION

## Conditions of Use

The Toolbox is free software; you can redistribute it and/or modify it under the terms of version 2 of [GNU General Public License](#) as published by the Free Software Foundation.

When *MIRtoolbox* is used for academic research, we would highly appreciate if scientific publications of works partly based on *MIRtoolbox* cite one of the following publications:

Olivier Lartillot, Petri Toivainen, “A Matlab Toolbox for Musical Feature Extraction From Audio”, [International Conference on Digital Audio Effects](#), Bordeaux, 2007.

Olivier Lartillot, Petri Toivainen, Tuomas Eerola, “A Matlab Toolbox for Music Information Retrieval”, in C. Preisach, H. Burkhardt, L. Schmidt-Thieme, R. Decker (Eds.), [Data Analysis, Machine Learning and Applications](#), Studies in Classification, Data Analysis, and Knowledge Organization, Springer-Verlag, 2008.

For commercial use of *MIRtoolbox*, please contact the authors.

## Please Register

Please register to the *MIRtoolbox* [announcement list](#).

This will allow us to estimate the number of users, and this will allow you in return to get informed on the new major releases (including critical bug fixes).

## Documentation and Support

The URL of *MIRtoolbox* website is [www.jyu.fi/music/coe/materials/mirtoolbox](http://www.jyu.fi/music/coe/materials/mirtoolbox)

### MIRTOOLBOX DISCUSSION LIST

A discussion list is also available:

- To subscribe, send an empty mail with ‘Subscribe’ as subject to [mirtoolbox-request@free-lists.org](mailto:mirtoolbox-request@free-lists.org)
- The archive is available [here](#).

## MIRTOOLBOX TWEETS

Get informed of the day-to-day advance of the project (bug reports, bug fixes, new features, new topics, etc.) by following [@mirtoolbox](#).

## TUTORIAL VIDEO

Video recordings of a tutorial given during SMC09 are available on [YouTube](#).

## Background

### ABOUT THE AUTHORS

Olivier Lartillot, Petri Toiviainen, Pasi Saari and Tuomas Eerola were members of the *Finnish Centre of Excellence in Interdisciplinary Music Research*, University of Jyväskylä, Finland.

The development of the toolbox has benefitted from productive collaborations with:

- partners of the *Brain Tuning* project (Marco Fabiani, Jose Fornari, Anders Friberg, Roberto Bresin, ...),
- colleagues from the *Finnish Centre of Excellence in Interdisciplinary Music Research* (Martin Hartmann, Vinoo Alluri, Rafael Ferrer, Marc Thompson, Anemone Van Zijl, Iballa Burunat, Birgitta Burger, Tiziana Quarto, ...),
- colleagues from the *Swiss Center for Affective Sciences* (Didier Grandjean, Klaus Scherer, Eduardo Coutinho, Donald Glowinski, Kim Eliard, Carolina Labbé, Johanna Wiebke Trost, Donato Cereghetti, ...),
- colleagues from the *RITMO Centre for Interdisciplinary Studies in Rhythm, Time and Motion*, University of Oslo, Norway (Kristian Nymoen, Anne Danielsen, Alexander Refsum Jensenius, Rolf Inge Godøy, Guilherme Schmidt Câmara, Merve Akca, ...),
- students of the MMT master program (University of Jyväskylä, Finland) and of the SMC master program (Aalborg University, Denmark),
- external collaborators: Stéphanie Weisser (Université Libre de Bruxelles), Patrik Juslin (Uppsala University), Jakob Abeßer (Fraunhofer IDMT), Thomas Wosch and associates (MEM, FHWS), Cyril Laurier and Emilia Gomez, (MTG-UPF),
- active users of the toolbox, participating in particular to the discussion list,
- participants of summer schools: SMC Summer School 2007 and 2014, ISSSM 2007, ISSSCCM 2009, USMIR 2010, ISSAS 2011 and 2012.

## TUNING THE BRAIN FOR MUSIC

MIRtoolbox has been initially developed within the context of a European Project called “*Tuning the Brain for Music*”, funded by the NEST (New and Emerging Science and Technology) program of the European Commission. The project, coordinated by Mari Tervaniemi from the Cognitive Brain Research Unit of the Department of Helsinki, is dedicated to the study of music and emotion, with collaboration between neurosciences, cognitive psychology and computer science. One particular question, studied in collaboration between the Music Cognition Team of the University of Jyväskylä and the Music Acoustics Group of the KTH in Stockholm, is related to the investigation of the relation between musical features and music-induced emotion. In particular, we would like to know which musical parameters can be related to the induction of particular emotion when playing or listening to music. For that purpose, we needed to extract a large set of musical features from large audio data-bases, in order to perform in a second step a statistical mapping between the diverse musical parameters and musical materials with listeners’ emotional ratings. This requires in particular a management of the interdependencies between the diverse features – in order to avoid having to recompute the same operations again and again – and also a control of the memory costs while analyzing the databases.

## MUSIC, MIND, TECHNOLOGY

Within the master degree, called *Music Mind Technology* (MMT) at the University of Jyväskylä, the *Music Information Retrieval* course, previously taught by Olivier Lartillot, Petri Toivainen and others, offers an overview of computer-based research for music analysis and in particular musical feature extraction. For the hands-on sessions, we wanted the student to be able to try by themselves the different computational approaches using *Matlab*. As many of them did not have much background in this programming environment, we decided to design a computational environment for musical feature extraction aimed at both expert and non-expert of *Matlab*.

## MIRtoolbox Objectives

Due to the context of development of this toolbox, we elaborated the following specifications:

### GENERAL FRAMEWORK

*MIRtoolbox* proposes a large set of musical feature extractors.

### MODULAR FRAMEWORK

*MIRtoolbox* is based on a set of building blocks that can be parametrized, reused, reordered, etc.

## SIMPLE AND ADAPTIVE SYNTAX

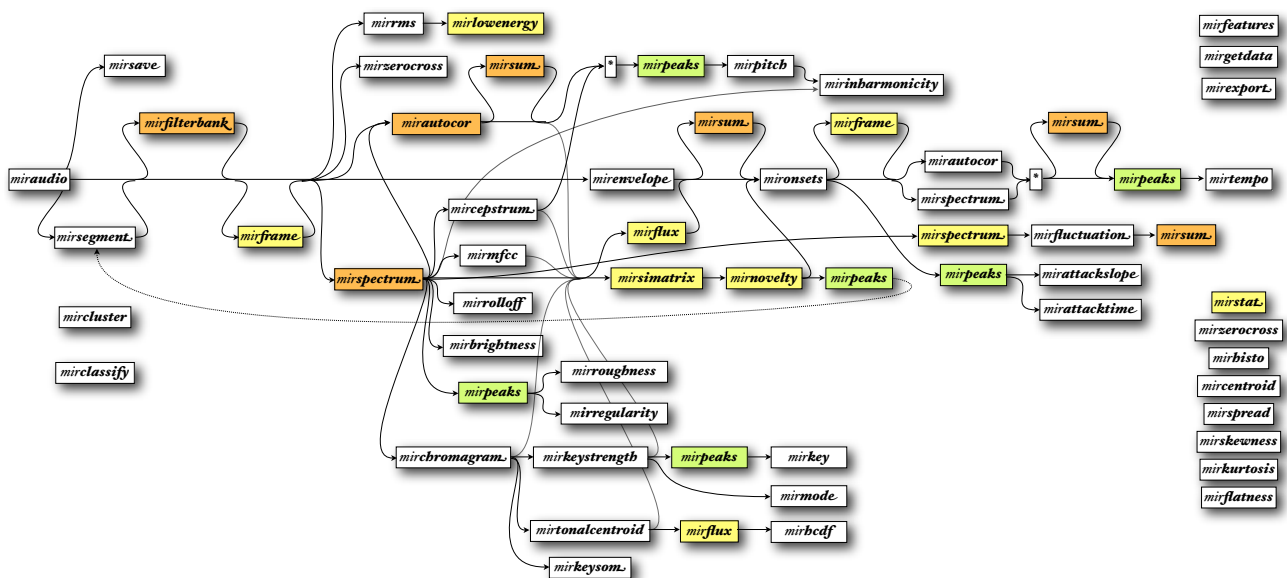
Users can focus on the general design, *MIRtoolbox* takes care of the underlying laborious tasks.

## FREE SOFTWARE, OPEN SOURCE

The idea is to propose to capitalize the expertise of the research community, and to offer it back to the community and the rest of us.

## *MIRtoolbox* Features

*MIRtoolbox* includes around 50 audio and music features extractors and statistical descriptors. A brief overview of most of the features can be seen in the following figure.



Synthetic overview of the features available in *MIRtoolbox* 1.2

## MIRtoolbox Reliances

### REQUIRED COMMERCIAL PRODUCTS

*MIRtoolbox* requires the **Matlab** environment, **version 7**, and does not work very well with previous versions of *Matlab*. This is due in particular to the fact *MIRtoolbox* relies on multi-dimensional arrays and multiple outputs, which seem to be features introduced by version 7.

*MIRtoolbox* also requires that the **Signal Processing Toolbox**, one of the optional sub-packages of *Matlab*, be properly installed. But actually, a certain number of operators can adapt to the absence of this toolbox, and can produce more or less reliable results. But for serious use of *MIRtoolbox*, we strongly recommend a proper installation of the *Signal Processing Toolbox*.

*MIRtoolbox* also requires that the **Statistics Toolbox (now called Statistics and Machine Learning Toolbox)**, one of the optional sub-packages of *Matlab*, be properly installed. This is needed for *mirflux* and for any other operator using that module.

### FREE SOFTWARES INCLUDED IN THE MIR TOOLBOX DISTRIBUTION

*MIRtoolbox* includes in its distribution several other freely available toolboxes, that are used for specific computations.

- The *Auditory Toolbox*, by Malcolm Slaney (1998), is used for Mel-band spectrum and MFCC computations, and Gammatone filterbank decomposition.
- The *Netlab* toolbox, by Ian Nabney (2002), where the routines for Gaussian Mixture Modeling (GMM) is used for classification (*mirclassify*).
- Finally, the SOM toolbox, by Esa Alhoniemi and colleagues (Vesanto, 1999), where only a routine for clustering based on *k*-means method is used, in the *mircluster* function.

### CODE INTEGRATED AS PART OF GPL PROJECT

*MIRtoolbox* license is based on GPL 2.0. As such, it can integrate codes from other GPL 2.0 projects, as long as their origins are explicitly stated.

- codes from the *Music Analysis Toolbox* by Elias Pampalk (2004), related to the computation of Terhardt outer ear modeling, Bark band decomposition and masking effects. (GPL 2.0)
- implementation of Earth Mover Distance written by Yossi Rubner and wrapped for Matlab by Simon Dixon.
- *openbdf* and *readbdf* script by T.S. Lorig to read BDF files, based on *openedf* and *readedf* by Alois Schloegl.

## CODE INTEGRATED WITH BSD LICENSE

- [\*mp3read\*](#) by Dan Ellis, which calls the [\*mpg123\*](#) decoder and the [\*mp3info\*](#) scanner
- [\*aiffread\*](#) by Kenneth Eaton
- [\*convolve2\*](#) by David Young

## Installation

To install *MIRtoolbox* in your *Matlab* environment, move the main *MIRtoolbox* folder to the location of your choice in your computer (for instance, in your *Matlab* "toolbox" folder, if you have administrative rights to modify it). Then open the "Set Path" environment available in *Matlab* File menu, click on "Add with Subfolders...", browse into the file hierarchy and select the main *MIRtoolbox* folder, then click "Open". You can then "Save" and "Close" the *Set Path* environment.

## UPDATE

If you replace an older version of *MIRtoolbox* with a new one, please update your *Matlab* path using the following command:

*rehash toolboxcache*

Update also the class structure of the toolbox, either by restarting *Matlab*, or by typing the following command:

*clear classes*

## MP3 READER FOR MAC OS X 64-BITS PLATFORM

The problem discussed in this section should not be an issue anymore if you use *Matlab* 2014a or more recent with *MIRtoolbox* 1.6 or more recent. Because in that case, *MIRtoolbox* uses the new audio reader introduced in *Matlab*. But beware that in this case **you need to make sure that you include the file extension of the audio file names you specify.**

If you are running *Matlab* on a *Mac OS X* 10.6 or beyond and with *Matlab* release 2009 or beyond, the binaries used for reading MP3 files (*mpg123* and *mp3info*) needs to be in 64-bits format (with the *mexmaci64* file extension). Unfortunately, it seems that the *mpg123.mexmaci64* and *mp3info.mexmaci64* executable we provided in the *MIRtoolbox* distribution cannot be used directly on other computers, so you need to install those binaries by yourselves on each separate computer by doing the following:

- Install Apple's Xcode:



- If you use Mac OS X 10.7, you can download it from free on the *Mac App Store*. After installing *Xcode*, it is advised to install the Command Line Tools must be installed. This is done from the Downloads section of *Xcode*'s preferences.
- If you use Mac OS X 10.6, you need to be (freely) registered as an Apple Developer. We suggest to download *Xcode* 3.2.6, as it is the latest free version available.
- Install *MacPorts*.
- Check that your *MacPorts* is up-to-date by executing in the Terminal:

```
sudo port -v selfupdate
```

(You need to authenticate as an administrative user.)

- Install *mpg123* and *mp3info* via *MacPorts* by executing in the Terminal:

```
sudo port install mpg123
```

```
sudo port install mp3info
```

- Once both installations are completed, you should obtain among others two Unix executable files called *mpg123* and *mp3info*, probably located at the address */opt/local/bin*.
- Create a copy of these files that you rename *mpg123.mexmaci64* and *mp3info.mexmaci64*, and place these two renamed files in a folder whose path is included in Matlab. You can for instance place them in your *MIRtoolbox* folder, which already contains Unix executable *mpg123.mexmaci* and *mp3info.mexmaci*, which correspond to the 32-bit platform. If there already exists files called *mpg123.mexmaci64* and *mp3info.mexmaci64*, you can replace those previous files with the new ones you compiled yourself.

## Help and demos

To get an overview of the functions available in the toolbox, type:

*help mirtoolbox*

A short documentation for each function is available using the same *help* command. For instance, type:

*help miraudio*

### D E M O S

Examples of use of the toolbox are shown in the *MIRToolboxDemos* folder:

- *mirdemo*
- *demo1basics*
- *demo2timbre*
- *demo3tempo*
- *demo4segmentation*
- *demo5export*
- *demo6curves*
- *demo7tonality*
- *demo8classification*
- *demo9retrieval*

## MIRtoolbox Interface

### BASIC SYNTAX

All functions are preceded by the *mir-* prefix in order to avoid conflicts with other Matlab functions. Each function is related to a particular data type: for instance, *miraudio* is related to the loading, transformation and display of audio waveform. An audio file, let's say a WAV file of name *mysong.wav*, can be loaded simply by writing the command:

*miraudio('mysong.wav')*

The extension of the file should be indicated explicitly. In recent versions of Matlab, the following syntax:

*miraudio('mysong')*

returns an error.

Operations and options to be applied are indicated by particular keywords, expressed as arguments of the functions. For instance, the waveform can be centered using the 'Center' keyword:

*miraudio('mysong.wav', 'Center')*

which is equivalent to any of these parameters:

*miraudio('mysong.wav', 'Center', 'yes')*

*miraudio('mysong.wav', 'Center', 'on')*

*miraudio('mysong.wav', 'Center', 1)*

whereas the opposite set of parameters

*miraudio('mysong.wav', 'Center', 'no')*

*miraudio('mysong.wav', 'Center', 'off')*

*miraudio('mysong.wav', 'Center', 0)*

are not necessary in the case of the ‘Center’ options as it is toggle off by default in *miraudio*.

It should be noted also that keywords are not case-sensitive:

*miraudio*(‘mysong.wav’, ‘center’, ‘YES’)

Other options accept numerical particular parameters. For instance, an audio waveform can be resampled to any sampling rate, which is indicated by a value in Hertz (Hz.) indicated after the ‘Sampling’ keyword. For instance, to resample at 11025 Hz., we just write:

*miraudio*(‘mysong.wav’, ‘**Sampling**’, 11025)

Finally the different options can be combined in one single command line:

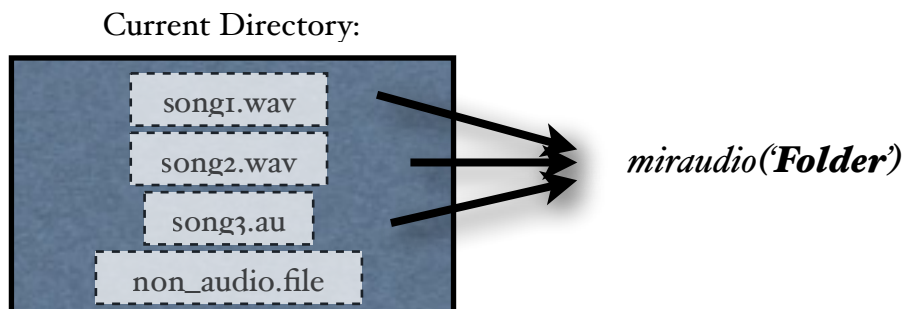
*miraudio*(‘mysong.wav’, ‘**Center**’, ‘**Sampling**’, 11025)

## B A T C H   A N A L Y S I S

- Folder of files can be analyzed in exactly the same way. For that, the file name, which was initially the first argument of the functions, can be replaced by the ‘**Folder**’ keyword. For instance, a folder of audio files can be loaded like this:

*miraudio*(‘**Folder**’)

Only audio files in audio formats recognized by Matlab’s in-house *audioread* function are taken into consideration, the other files are simply ignored:



*Automatic analysis of a batch of audio files using the ‘Folder’ keyword*

- Subfolders can be analyzed recursively as well, using the ‘**Folders**’ keyword:

*miraudio*(‘**Folders**’)

- Alternatively, the list of audio files (with their respective path) can be stored in successive lines of a TXT file, whose name (and path) can be given as input to *miraudio*:

*`miraudio('myfilenames.txt')`*

- As another alternative, the list of audio files (with address relative to the current directory) can be given in a cell array as first input to *miraudio*, for instance:

*`miraudio({'song1.wav', 'song2.au', 'song3.mp3'})`*

## OUTPUT FORMAT

After entering one command, such as

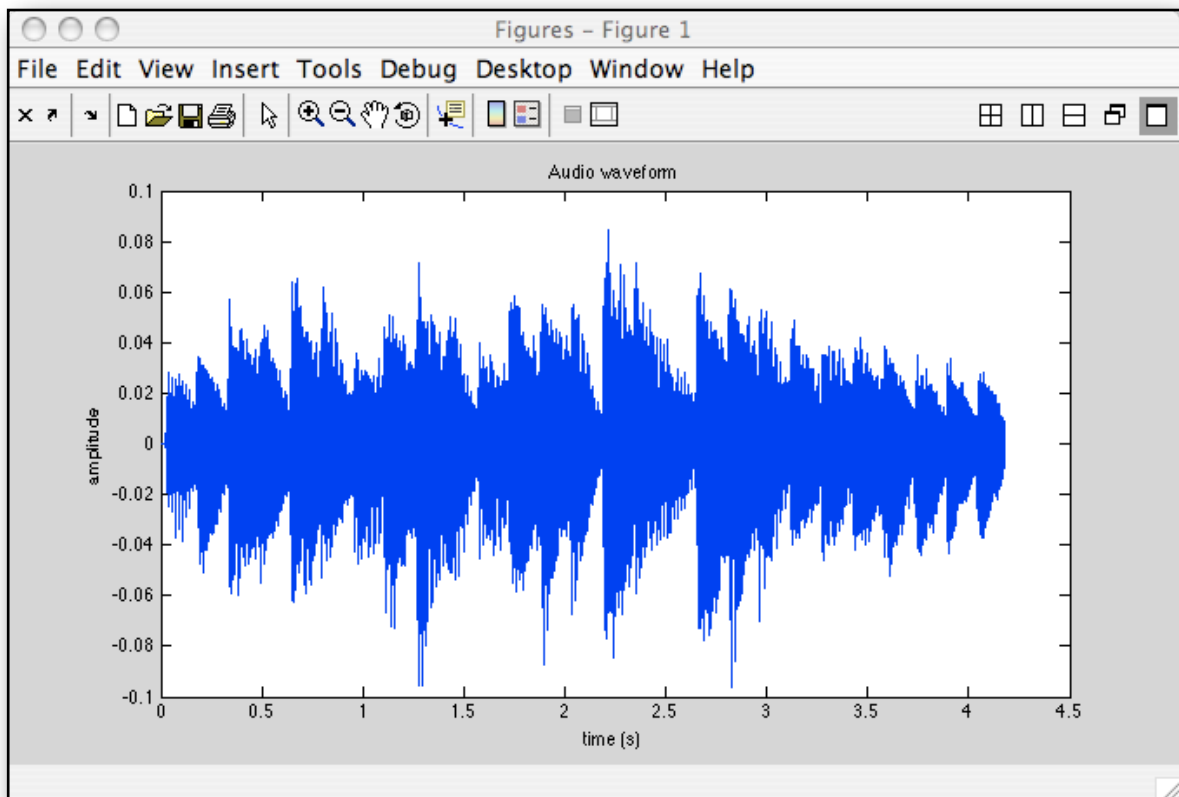
*`miraudio('mysong.wav')`*

the computation is carried out, and when it is completed, a text is written in the Command Window:

ans is the Audio waveform related to file mysong.wav, of sampling rate 44100 Hz.

Its content is displayed in Figure 1.

And a graphical representation of the result is displayed in a figure:



*Display of a miraudio object*

The display of the figures and the messages can be avoided, if necessary, by adding a semi-colon at the end of the command:

***miraudio('mysong.wav');***

The actual output is stored in an object, hidden by default to the users, which contains all the information related to the data, such as the numerical values of the waveform amplitudes, the temporal dates of the bins, the sampling rate, the name of the file, etc. In this way we avoid the traditional interface in Matlab, not quite user-friendly in this respect, where results are directly displayed in the Command Window by a huge list of numbers.

It is not possible to display MIRtoolbox results in the Matlab Variable Editor. If you try visualizing a MIRtoolbox variable listed in your Workspace window, for instance the audio waveform in the previous example, you get the following text in the Variable Editor:

```
val is the Audio waveform related to file mysong.wav, of sampling  
rate 44100 Hz.
```

```
To display its content in a figure, evaluate this variable directly  
in the Command Window.
```

## MULTIPLE FILE OUTPUT

If we now analyze a folder of file:

***miraudio('Folder')***

the results related to each audio file is displayed in a different figure, and messages such as the following ones are displayed in the Command Window:

```
ans(1) is the Audio waveform related to file song1.wav, of sampling  
rate 44100 Hz.
```

```
Its content is displayed in Figure 1.
```

```
ans(2) is the Audio waveform related to file song2.wav, of sampling  
rate 22050 Hz.
```

```
Its content is displayed in Figure 2.
```

```
ans(3) is the Audio waveform related to file song3.au, of sampling  
rate 11025 Hz.
```

```
Its content is displayed in Figure 3.
```

and so on.

And the actual output is stored in one single object, that contains the information related to all the different audio files.

## THREADING OF DATA FLOW

The result of one operation can be used for subsequent operations. For that purpose, it is better to store each result in a variable. For instance, the audio waveform(s) can be stored in one variable *a*:

$$\mathbf{a} = \text{miraudio}(\textit{mysong.wav});$$

Then the spectrum, for instance, related to the audio waveform can be computed by calling the function *mirspectrum* using simply the *a* variable as argument:

$$s = \text{mirspectrum}(\mathbf{a})$$

In this way, all the information necessary for the computation of the spectrum can be retrieved from the hidden object, associated to the variable *a*, that contains the complex encapsulated data.

Alternatively, the spectrum can be directly computed from a given audio file by indicating the file name as argument of the *mirspectrum* function:

$$s = \text{mirspectrum}(\textit{mysong.wav})$$

This second syntax, more compact, is generally recommended, because it avoids the decomposition of the computation in several steps (*a*, then *s*, etc.), which might cause significant problems for long audio files or for folder of files. We will see in section 5.3 how to devise more subtle datacharts that take into account memory management problems in a more efficient way.

## SUCCESSIVE OPERATIONS ON ONE SAME DATA FORMAT

When some data has been computed on a given format, let's say an audio waveform using the *miraudio* function:

$$\mathbf{a} = \text{miraudio}(\textit{mysong.wav});$$

it is possible to apply options related to that format in successive step. For instance, we can center the audio waveform in a second step:

$$a = \text{miraudio}(\mathbf{a}, \textit{Center});$$

which could more efficiently be written in one single line:

```
a = miraudio('my song.wav', 'Center');
```

## GETTING NUMERICAL DATA

The numerical data encapsulated in the output objects can be recuperated if necessary. In particular, the main numerical data (such as the amplitudes of the audio waveform) are obtained using the *mirgetdata* command:

```
mirgetdata(a)
```

the other related informations are obtained using the generic *get* method. For instance, the sampling rate of the waveform *a* is obtained using the command:

```
get(a, 'Sampling')
```

More detailed description of these functions will be described in section 5, dedicated to advanced uses of *MIRtoolbox*.

## SPEEDING UP MIRTOOLBOX

If you notice that *MIRtoolbox* takes too much time decomposing each files into a large sets of chunks, with a large displays of messages of the form “Chunk 1/20...” for instance, you might try to increase the chunk size by using *mirchunklim* (cf. section 5.5).



## 2. BASIC OPERATORS

*MIRtoolbox* basic operators concern the management of audio waveforms (*miraudio*, *mirsave*), frame-based analysis (*mirframe*, *mirflux*), periodicity estimation (*mirautocor*, *mirspectrum*, *mircepstrum*), operations related more or less to auditory modeling (*mirenvelope*, *mirfilterbank*), peak picking (*mirpeaks*) and sonification of the results (*mirplay*).

### *miraudio*

#### AUDIO WAVEFORM

As explained previously, this operator basically loads audio files, displays and performs operations on the waveform.

#### ACCEPTED INPUT FORMATS

- **file name:**

- if you use Matlab 2014a or more recent, **and if you make sure that you include the file extension when specifying the file name**, then *miraudio* uses Matlab's *audioread*, which accepts a large range of audio file format. They are listed in this page: [http://se.mathworks.com/help/matlab/ref/audioread.html#inputarg\\_filename](http://se.mathworks.com/help/matlab/ref/audioread.html#inputarg_filename)
- ~~in the other cases, the accepted file format are WAV, MP3, AIFF and AU formats, as the loading operations are based on the *Matlab wavread* and *auread* functions, on Dan Ellis' *mp3read* and on Kenneth Eaton's *aiffread*. But beware that the compiled files included in *MIRtoolbox* distribution (for instance *mp3read*) do not necessarily correspond to your platform. So you might get error when, for instance, trying to read a MP3 file. (Not working on recent Matlab versions.)~~

- ***miraudio* object:** for further transformations.

- ***Matlab* array:** It is possible to import an audio waveform encoded into a *Matlab* column vector, by using the following syntax:

*miraudio*(*v*, *sr*)

where *v* is a column vector and *sr* is the sampling rate of the signal, in Hz. The default value for *sr* is 44100 Hz.

## TRANSFORMATION OPTIONS

- *miraudio*(..., '**Mono**', *o*) does not perform the default summing of channels into one single mono track, but instead stores each channel of the initial sound file separately.
- *miraudio*(..., '**Center**') centers the waveform.
- *miraudio*(..., '**Sampling**', *r*) resamples at sampling rate *r* (in Hz). It uses the *resample* function from *Signal Processing Toolbox*.
- *miraudio*(..., '**Normal**', *n*) normalises according to a chosen heuristic:
  - *n* = '**RMS**': normalises with respect to RMS energy (cf. *mirrms*). (Default option when '**Normal**' is called.)
  - *n* = '**Max**': normalises with respect to the maximum of the absolute values.
- *miraudio*(..., '**FWR**') performs a full wave rectification, by flipping the negative values into positive values (absolute values).
- *miraudio*(..., '**Frame**', *w*, *wu*, *h*, *hu*) decomposes into frames. Cf. *mirframe* for an explanation of the arguments (units can be omitted here as well). Default parameters: same as in *mirframe*, i.e., 50 ms and half-overlapping.
- *miraudio*(..., '**SmoothBorder**', *l*) avoids clicks at the beginning and end of the signals by adding a Hann(ing) half-windowing just at the very beginning and end of the signal. The half-window is of duration *l* in milliseconds. Default parameter for *l* = 1 ms.

## EXTRACTION OPTIONS

- *miraudio*(..., '**Extract**', *t1*, *t2*, *u*, *f*) extracts the signal between the dates *t1* and *t2*, expressed in the unit *u*.
  - Possible units *u* = '*s*' (seconds, by default) or *u* = '*sp*' (sample index, starting from 1).
  - The additional optional argument *f* indicates the referential origin of the temporal positions. Possible values for *f*:
    - '*Start*' (by default),
    - '*Middle*' (of the sequence),
    - '*End*' of the sequence.

When using '*Middle*' or '*End*', negative values for *t1* or *t2* indicate values before the middle or the end of the audio sequence. For instance: *miraudio*(..., '**Extract**', -1, +1, '*Middle*') extracts one second before and after the middle of the audio file.

- Alternative keyword: ***Excerpt***.
- *miraudio(..., 'Trim')* trims the pseudo-silence beginning and end off the audio file.
  - *miraudio(..., 'TrimThreshold', t)* specifies the trimming threshold  $t$ . Silent frames are frames with RMS energy below  $t$  times the medium RMS of the whole audio file. Default value:  $t = 0.06$ .
  - Instead of 'Trim', '**TrimStart**' only trims the beginning of the audio file, whereas '**TrimEnd**' only trims the end.
- *miraudio(..., 'Channel', c)* or *miraudio(..., 'Channels', c)* selects the channels indicated by the (array of) integer(s)  $c$ .

## LABELING OPTION

*miraudio(..., 'Label', lb)* labels the audio signals following the name of their respective audio files.  $lb$  is one number, or an array of numbers, and the audio signals are labelled using the substring of their respective file name of index  $lb$ . If  $lb = 0$ , the audio signal(s) are labelled using the whole file name.

<i>miraudio('Folder', 'Label', lb)</i>	<i>song1g.wav</i>	<i>song2g.wav</i>	<i>song3b.au</i>
$lb = 6$	'g'	'g'	'b'
$lb = [5 \ 6]$	'1g'	'2g'	'3b'
$lb = \{'good', 'bad'\}$	'good'	'bad'	'good'

### Example of labelling of a folder of audio files

The labeling is used for classification purposes (cf. *mirclassify* and *mirexport*).

## SUMMATION

Audio signals can be superposed using the basic *Matlab* summation operators (+). For instance let's say we have two sequences:

$a1 = \text{miraudio('melody.wav')};$

$a2 = \text{miraudio('accompaniment.wav')};$

Then the two sequences can be superposed using the command:

$$a = a1+a2$$

When superposing *miraudio* objects, the longest audio are no more truncated, but on the contrary the shortest one are prolonged by silence. When audio have different sampling rates, all are converted to the highest one.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- ***Time***: the temporal positions of samples (same as *Pos*),
- ***Centered***: whether the waveform has been centered (1) or not (0),
- ***NBits***: the number of bits used to code each sample,
- ***Label***: the label associated to each audio file.

## *mirframe* vs. 'Frame'

### FRAME DECOMPOSITION

The analysis of a whole temporal signal (such as an audio waveform in particular) leads to a global description of the average value of the feature under study. In order to take into account the dynamic evolution of the feature, the analysis has to be carried out on a short-term window that moves chronologically along the temporal signal. Each position of the window is called a frame.

### FLOWCHART INTERCONNECTIONS

*mirframe* accepts as input any temporal object:

- an audio waveform *miraudio*,
- **file name** or the '**Folder**' keyword,
- an envelope *mirenvelope*,
- the **temporal** evolution of a scalar data, such as fluxes in particular (*mirflux*),
- in particular, event detection curves (*miरेvents*) can be decomposed into frames as well.

### SYNTAX

The frame decomposition can be performed using the *mirframe* command. The frames can be specified as follows:

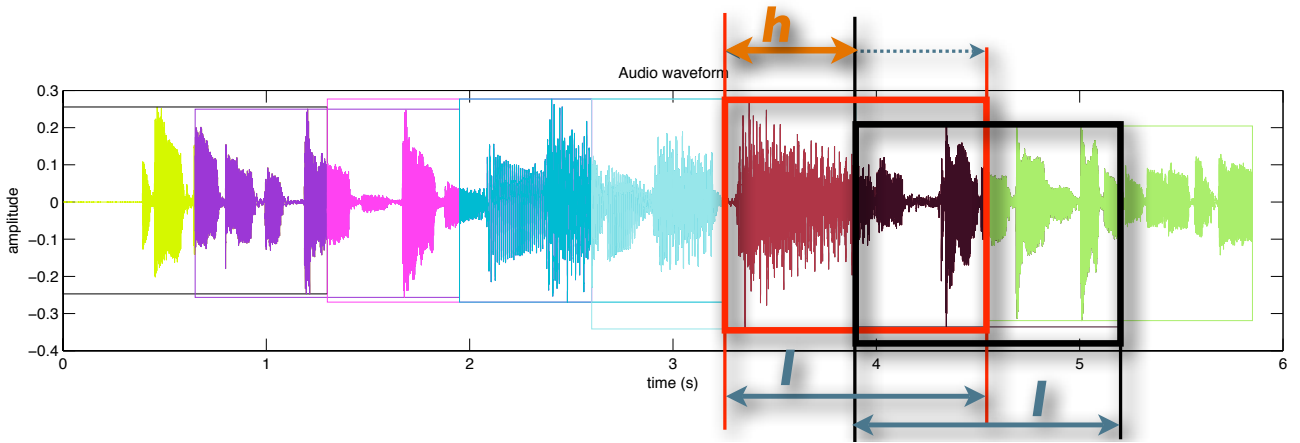
*mirframe*(*x*,..., '**Length**', *w*, *wu*):

- *w* is the length of the window in seconds (default: .05 seconds);
- *u* is the unit, either
  - '**s**' (seconds, default unit),
  - or '**sp**' (number of samples).

*mirframe*(*x*, '...', '**Hop**', *b*, *bu*):

- *b* is the hop factor, or distance between successive frames (default: half overlapping: each frame begins at the middle of the previous frame)
- *u* is the unit, either

- **'/l'** (ratio with respect to the frame length, default unit)
- **'%'** (ratio as percentage)
- **'s'** (seconds)
- **'sp'** (number of samples)
- or **'Hz'** (hertz), i.e., number of frames per second: the exactness of the frame rate is ensured and may cause a slight fluctuation of the elementary hop distances.



*Frame decomposition of an audio waveform, with frame length  $l$  and hop factor  $h$  (represented here, following the default unit, as a ratio with respect to the frame length).*

These arguments can also be written as follows (where units can be omitted):

$$\text{mirframe}(x, w, wu, h, hu)$$

## CHAINING OF OPERATIONS

Suppose we load an audio file:

$$a = \text{miraudio}(\text{'mysong'})$$

then we decompose into frames

$$f = \text{mirframe}(a)$$

then we can perform any computation on each of the successive frame easily. For instance, the computation of the spectrum in each frame (or spectrogram), can be written as:

$$s = \text{mirspectrum}(f)$$

## THE ‘FRAME’ OPTION

The two first previous commands can be condensed into one line, using the ‘Frame’ option.

$$f = \text{miraudio}(\text{'mysong'}, \text{'Frame'})$$

and the three commands can be condensed into one line also using the ‘Frame’ option.

$$s = \text{mirspectrum}(\text{'mysong'}, \text{'Frame'})$$

The frame specifications can be expressed in the following way:

$$\text{mirspectrum}(\dots, \text{'Frame'}, l, s, b, /I)$$

It is also possible to specify just the hop factor using the following syntax:

$$\text{mirspectrum}(\dots, \text{'Frame'}, \text{'Hop'}, h, /I)$$

This ‘Frame’ option is available to most operators. Each operator uses specific default values for the ‘Frame’ parameters. Each operator can perform the frame decomposition where it is most suitable. For instance, as can be seen in *mirevents* feature map, the ‘Frame’ option related to the *mirevents* operator will lead to a frame decomposition after the actual computation of the event detection curve (produced by *mirevents*).

## ACCESSIBLE OUTPUT

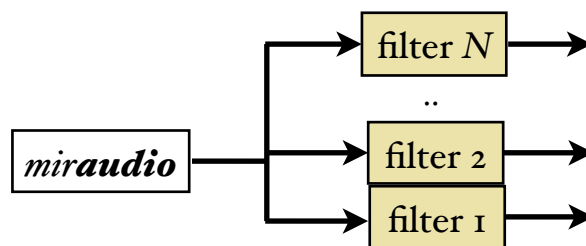
cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **‘FramePos’**: the starting and ending temporal positions of each successive frame, stored in the same way as for ‘Data’ (cf. §5.2),
- **‘FrameRate’**: the frame rate in Hertz, i.e. the number of frames computed for each second,
- **‘Framed’**: whether the data has been decomposed into frames or not.

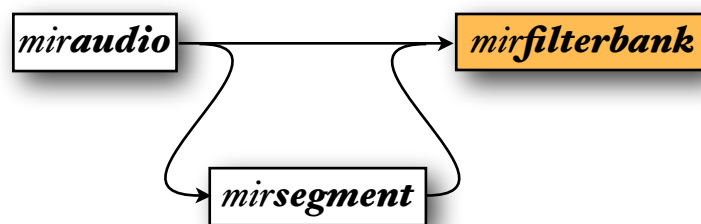
## *mirfilterbank*

### FILTERBANK DECOMPOSITION

It is often interesting to decompose the audio signal into a series of audio signals of different frequency register, from low frequency channels to high frequency channels. This enables thus to study each of these channels separately. The decomposition is performed by a bank of filters, each one selecting a particular range of frequency values. This transformation models an actual process of human perception, corresponding to the distribution of frequencies into critical bands in the cochlea.



### FLOWCHART INTERCONNECTIONS



*mirfilterbank* accepts as input data type either:

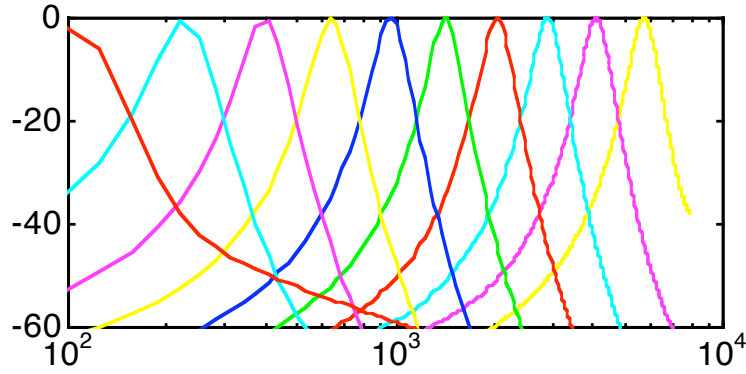
- ***miraudio*** objects, where the audio waveform can be segmented (using *mirsegment*),
- **file name** or the **'Folder'** keyword.

### FILTERBANK SELECTION

Two basic types of filterbanks are proposed in *MIRtoolbox*:

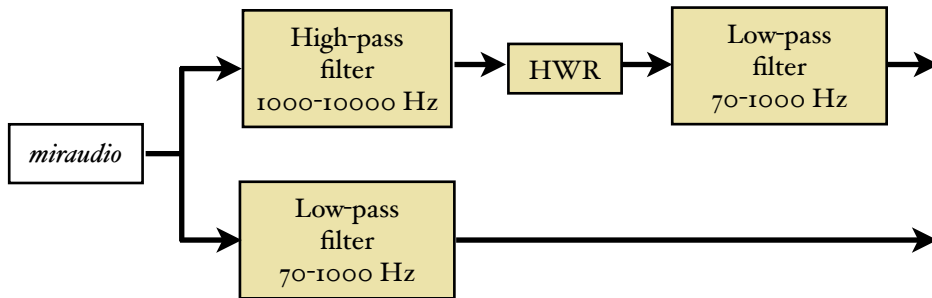
- *mirfilterbank*(..., '**Gammatone**') carries out a Gammatone filterbank decomposition (Patterson et al, 1992). It is known to simulate well the response of the basilar membrane. It is based on a Equivalent Rectangular Bandwidth (ERB) filterbank, meaning that the width of each band is determined by a particular psychoacoustical law. For Gammatone filterbanks, *mirfilterbank* calls the *Auditory Toolbox* routines *MakeERBFilters* and *ERBfilterbank*. This is the default choice when calling *mirfilterbank*.





*Ten ERB filters between 100 and 8000Hz (Slaney, 1998)*

- `mirfilterbank(..., 'Lowest', f)` indicates the lowest frequency  $f$ , in Hz. Default value: 50 Hz.
- `mirfilterbank(..., '2Channels')` performs a computational simplification of the filterbank using just two channels, one for low-frequencies, below 1000 Hz, and one for high-frequencies, over 1000 Hz (Tolonen and Karjalainen, 2000). On the high-frequency channel is performed an envelope extraction using a half-wave rectification and the same low-pass filter used for the low-frequency channel. This filterbank is mainly used for multi-pitch extraction (cf. *mirpitch*).



*Diagram of the two-channel filterbank proposed in (Tolonen and Karjalainen, 2000)*

For these general type of filterbanks are chosen, further options are available:

- `mirfilterbank(..., 'NbChannels', N)` specifies the number of channels in the bank. By default:  $N = 10$ . This option is useless for '2Channels'.
- `mirfilterbank(..., 'Channel', c)` – or `mirfilterbank(..., 'Channels', c)` – only output the channels whose ranks are indicated in the array  $c$ . (default:  $c = (1:N)$ )

## MANUAL SPECIFICATIONS

`mirfilterbank(..., 'Manual', f)` specifies an array of non-overlapping low-pass, band-pass and high-pass elliptic filters (Scheirer, 1998). The array of cut-off frequencies  $f$  as to be specified as next parameter.

- If this array of frequencies begins with `-Inf`, the first filter is low-pass.
- If this array of frequencies ends with `Inf`, the last filter is high-pass.

`mirfilterbank(..., 'Order', o)` specifies the order of the filters. The default is set to  $o = 4$  (Scheirer, 1998)

`mirfilterbank(..., 'Hop', h)` specifies the degree of spectral overlapping between successive channels.

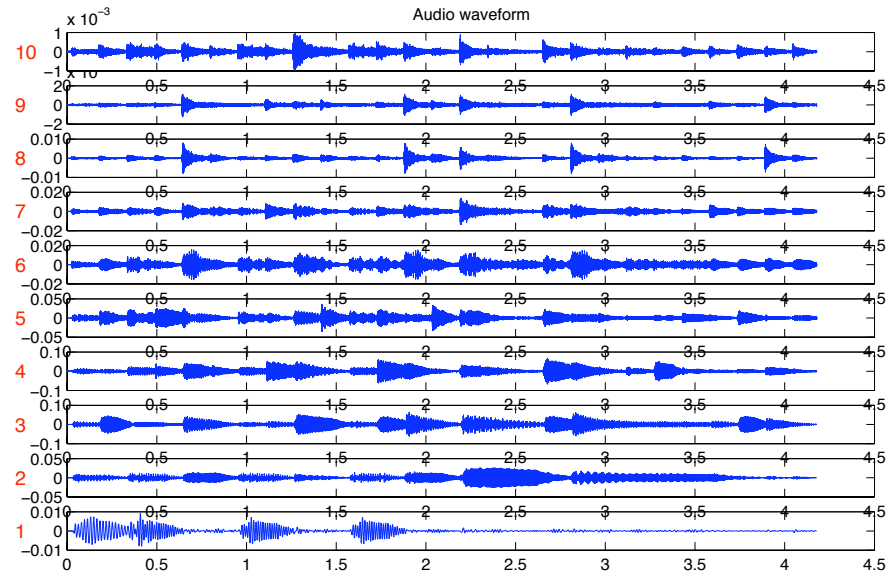
- If  $h = 1$  (default value), the filters are non-overlapping.
- If  $h = 2$ , the filters are half-overlapping.
- If  $h = 3$ , the spectral hop factor between successive filters is a third of the whole frequency region, etc.

## PRESELECTED FILTERBANKS

`mirfilterbank(..., p)` specifies predefined filterbanks, all implemented using elliptic filters, by default of order 4:

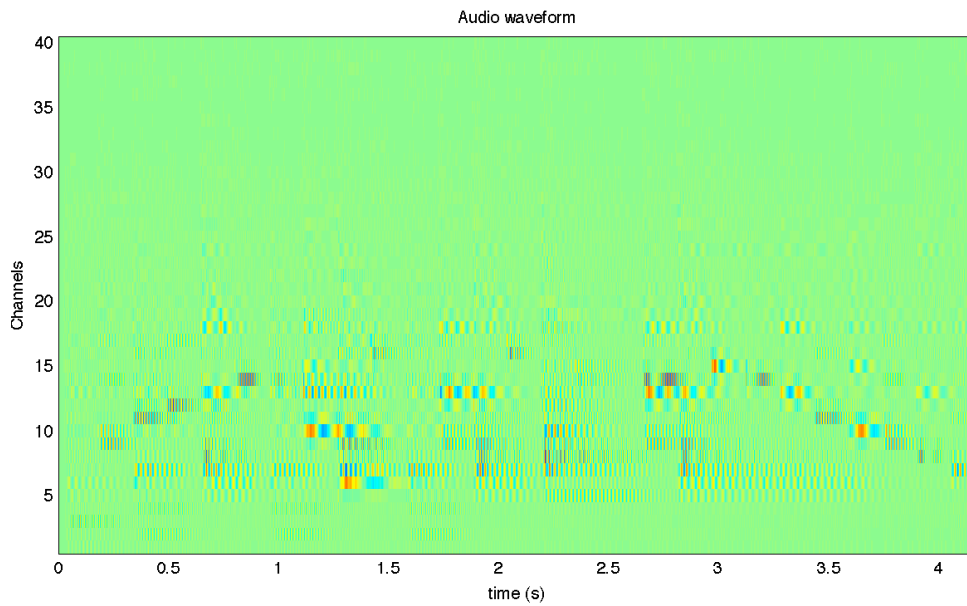
- $p = \text{'Mel'}$ : Mel scale (cf. `mirspectrum(..., 'Mel')`).
- $p = \text{'Bark'}$ : Bark scale (cf. `mirspectrum(..., 'Bark')`).
- $p = \text{'Scheirer'}$  proposed in (Scheirer, 1998) corresponds to `'Manual', [-Inf 200 400 800 1600 3200 Inf]`
- $p = \text{'Klapuri'}$  proposed in (Klapuri, 1999) corresponds to `'Manual', 44*[2.^(([ 0:2, (9+(0:17))/3 ]))]`

## EXAMPLE



*mirfilterbank('ragtime.wav.wav')*

If the number of channels exceeds 20, the audio waveform decomposition is represented as a single image bitmap, where each line of pixel represents each successive channel:



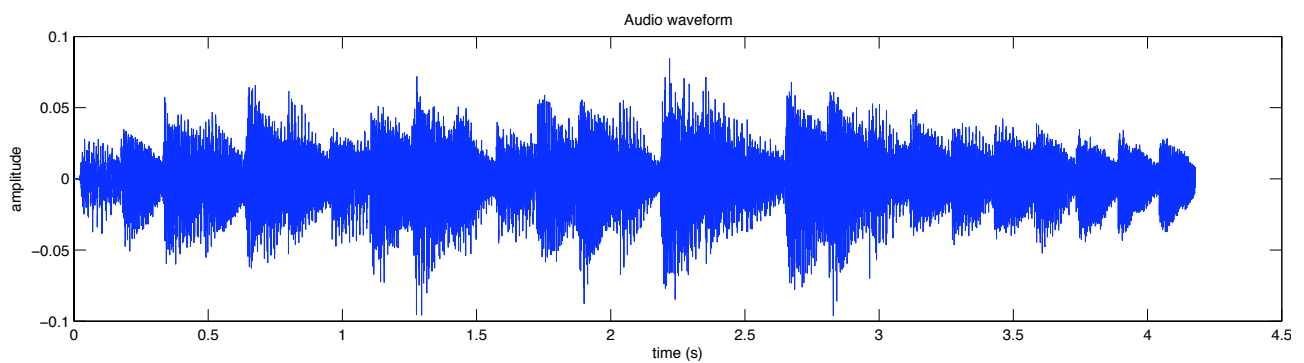
*mirfilterbank('ragtime.wav', 'NbChannels', 40)*

## *mirenvelope*

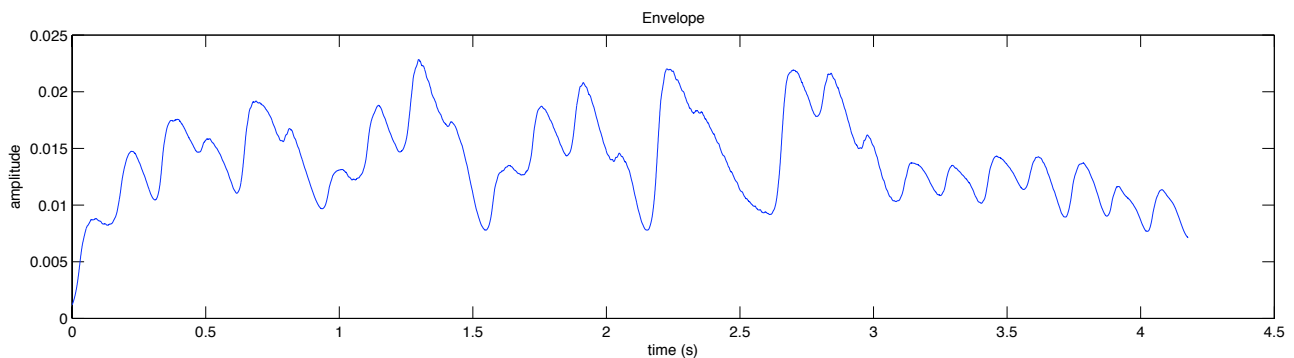
### AMPLITUDE ENVELOPE

From an audio waveform can be computed the envelope, which shows the global outer shape of the signal. It is particularly useful in order to show the long term evolution of the signal, and has application in particular to the detection of musical events such as notes.

Here is an example of audio file with its envelope:

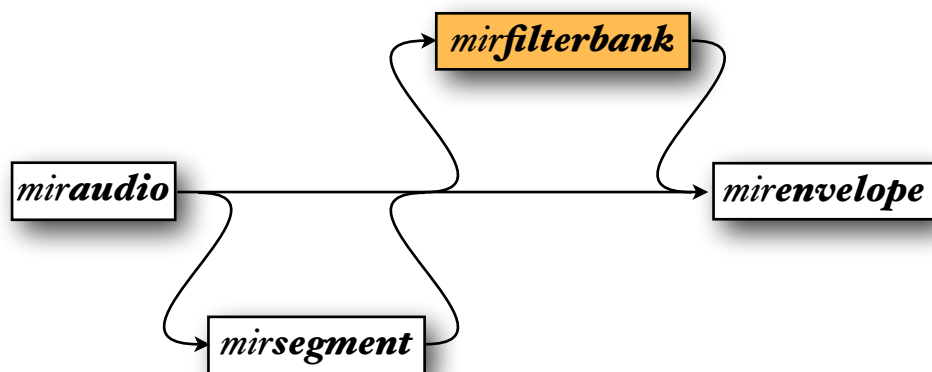


*Audio waveform of ragtime excerpt*



*Corresponding envelope of the ragtime excerpt*

### FLOWCHART INTERCONNECTIONS



*mirenvelope* accepts as input data type either:

- ***miraudio*** objects, where the audio waveform can be segmented (using ***mirsegment***) and/or decomposed into channels (using ***mirfilterbank***),
- **file name** or the **'Folder'** keyword,
- any scalar object (i.e., where there is one numerical value associated to each successive frame, such as *mirflux*, *mirnovelty*, etc.): in this case, the *mirscalar* object is simply converted into a *mirenvelope* object. The advantages of this operation is that the resulting *mirenvelope* can be further decomposed into frames, which would not have been possible using the *mirscalar* object as it is already decomposed into frames.

## PARAMETERS SPECIFICATION

The envelope extraction is based on two alternate strategies: either based on a filtering of the signal (**'Filter'** option, used by default), or on a decomposition into frames via a spectrogram computation (**'Spectro'** option). Each of these strategies accepts particular options:

- *mirenvelope*(..., **'Filter'**) extract the envelope through a filtering of the signal. (Default method.)
  - First the signal can be converted from the real domain to the complex domain using a Hilbert transform. In this way the envelope is estimated in a three-dimensional space defined by the product of the complex domain and the temporal axis. Indeed in this representation the signal looks like a “spring” of varying width, and the envelope would correspond to that varying width. In the real domain, on the other hand, the constant crossing of the signal with the zero axis may sometime give erroneous results.

An Hilbert transform can be performed in *mirenvelope*, based on the *Matlab* function *hilbert*. In order to toggle on the Hilbert transform, the following keyword should be added:

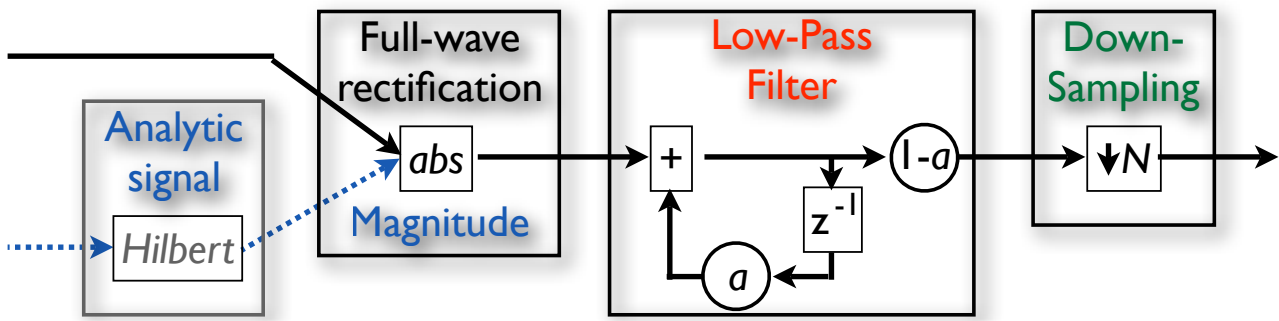
*mirenvelope*(..., **'Hilbert'**)

Beware however that, although sometimes the use of the Hilbert transform seems to improve somewhat the results, and might in particular show clearer burst of energy, we noticed some problematic behavior, in particular at the beginning and the end of the signal, and after some particular bursts of energy. This becomes all the more problematic when chunk decompositions are used (cf. §5.3), since the continuity between chunk cannot be ensured any more. For that reason, since version 1.1 of *MIRtoolbox*, the use of Hilbert transform is toggled off by default.

If the signal is in the real domain, the next step consists in a full-wave rectification, reflecting all the negative lobes of the signal into the positive domain, leading to a series of

positive half-wave lobes. The further smoothing of the signal (in the next step) will lead to an estimation of the envelope. If on the contrary the signal is in the complex domain, a direct estimation of the envelope can be obtained by computing the modulus, i.e., the width of the “string”. These two operations, either from the real or the complex domains, although apparently different, relate to the same *Matlab* command *abs*.

- *mirenvelope*(..., '**PreDecim**', *N*) down-samples by a factor  $N > 1$ , where *N* is an integer, before the low-pass filtering (Klapuri, 1999). Default value:  $N = 1$ , corresponding to no down-sampling.
- The next step consists in a low-pass filter that retains from the signal only the long-term evolution, by removing all the more rapid oscillations. This is performed through a filtering of the signal. Three types of filters are available:
  - *mirenvelope*(..., '**FilterType**', '**IIR**') extracts the envelope using an auto-regressive filter of infinite impulse response (IIR). This is the default method.



*Detail of the 'IIR' envelope extraction process*

The range of frequencies to be filtered can be controlled by selecting a proper value for the *a* parameter. Another way of expressing this parameter is by considering its time constant. If we feed the filter with a step function (i.e. 0 before time 0, and 1 after time 0), the time constant will correspond to the time it will take for the output to reach 63 % of the input. Hence higher time constant means smoother filtering. The default time constant is set to .02 seconds and can be changed using the option:

*mirenvelope*(..., '**Tau**', *t*)

- *mirenvelope*(..., '**FilterType**', '**HalfHann**') extracts the envelope using a half-Hanning (raised cosine) filter.
- *mirenvelope*(..., '**FilterType**', '**Butter**') extract the envelope using a Butterworth filter.
  - *mirenvelope*(..., '**CutOff**', *c*) controls the cut-off frequency, by default set to 37 Hz (Nymoen et al., 2017).

- Once the signal has been smoothed, as there is a lot of redundancy between the successive samples, the signal can be down-sampled. The default parameter related to down-sampling is the down-sampling rate  $N$ , i.e. the *integer* ratio between the old and the new sampling rate.  $N$  is set by default to 16, and can be changed using the option:

*mirenvelope(..., 'PostDecim', N)*

Alternatively, any sampling rate  $r$  (in Hz) can be specified using the post-processing option 'Sampling'.

- *mirenvelope(..., 'Trim')*: trims the initial ascending phase of the curves related to the transitory state.
- *mirenvelope(..., 'PreSilence')* adds further silence at the beginning of the audio sequence.
- *mirenvelope(..., 'PostSilence')* adds further silence at the end of the audio sequence.
- *mirenvelope(..., 'Spectro')* extracts the envelope through the computation of a power spectrogram, with frame size 100 ms, hop factor 10% and the use of Hanning windowing:

*mirspectrum(..., 'Frame', .1, 's', .1, '/I', 'Window', 'hanning', b, 'Power')*

Each frequency (or band) is considered as a specific channel of a multichannel representation. As such, channels can be summed back using *mirsum* (cf. *mirsum*).

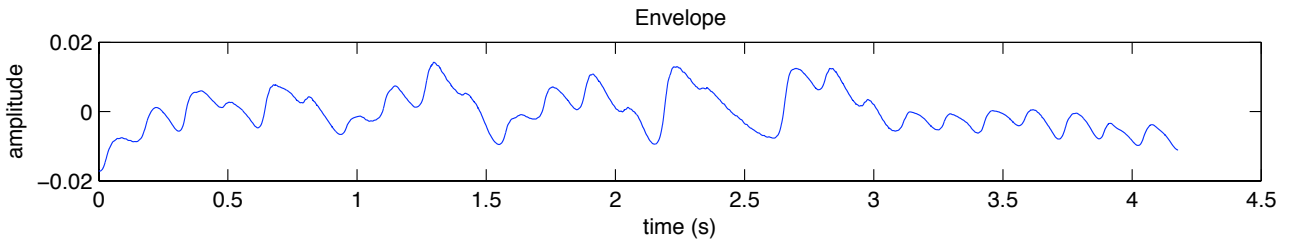
- *mirenvelope(..., b)* specifies whether the frequency range is further decomposed into bands (cf. *mirspectrum*). Possible values:
  - $b = 'Freq'$ : no band decomposition (default value),
  - $b = 'Mel'$ : Mel-band decomposition,
  - $b = 'Bark'$ : Bark-band decomposition,
  - $b = 'Cents'$ : decompositions into cents.
- *mirenvelope(..., 'Frame', ...)* modifies the default frame configuration.
- *mirenvelope(..., 'UpSample', N)* upsamples by a factor  $N > 1$ , where  $N$  is an integer. Default value if 'UpSample' called:  $N = 2$
- *mirenvelope(..., 'Complex')* toggles on the 'Complex' method for the spectral flux computation (cf. *mirflux*).
- *mirenvelope(..., 'PowerSpectrum', 0)* turns off the computation of the power of the spectrum.

- *mirenvelope*(..., '**Terhardt**') toggles on the '**Terhardt**' operation (cf. *mirspectrum*).
- *mirenvelope*(..., '**TimeSmooth**', *n*) toggles on and controls the '**TimeSmooth**' operation. (cf. *mirspectrum*).
- *mirenvelope*(..., '**PreSilence**') adds further silent frames before the actual start of the sequence.
- *mirenvelope*(..., '**PostSilence**') adds further silent frames after the actual end of the audio sequence.

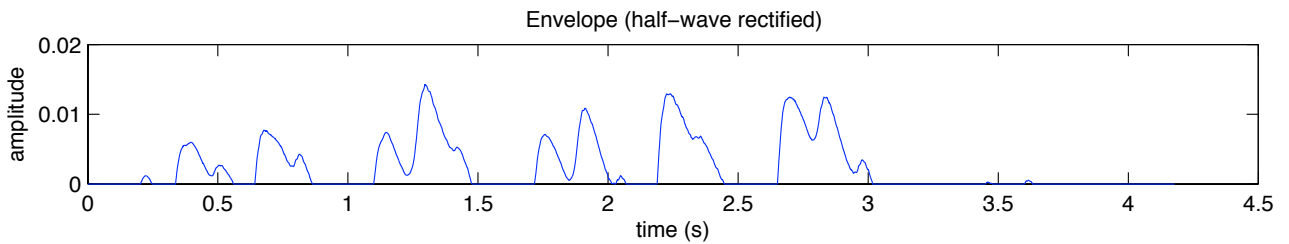
## POST-PROCESSING OPTIONS

Different operations can be performed on the envelope curve:

- *mirenvelope*(..., '**Sampling**', *r*) resamples to rate *r* (in Hz). '**PostDecim**' and '**Sampling**' options cannot therefore be combined.
- *mirenvelope*(..., '**Halfwave**') performs a half-wave rectification on the envelope.
- *mirenvelope*(..., '**Center**') centers the extracted envelope.



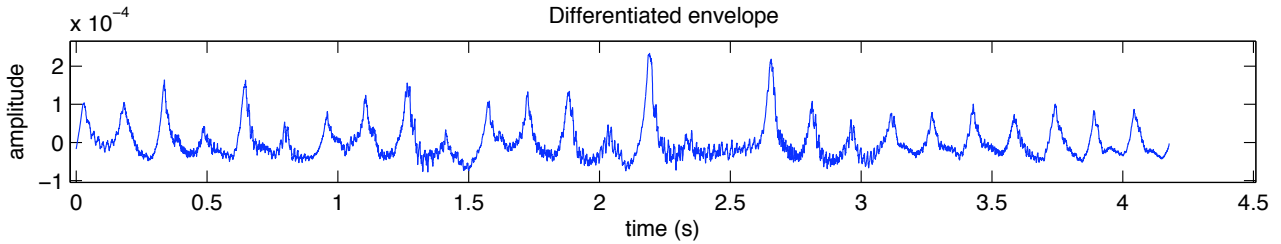
- *mirenvelope*(..., '**HalfwaveCenter**') performs a half-wave rectification on the centered envelope.



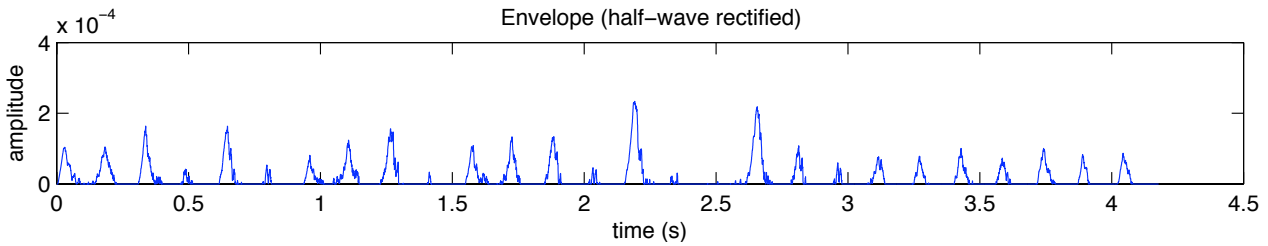
- *mirenvelope*(..., '**Log**') computes the common logarithm (base 10) of the envelope.
  - *mirenvelope*(..., '**MinLog**', *ml*) selects the data in the logarithmic range  $[-ml \text{ dB}, 0 \text{ dB}]$ , where 0 dB corresponds to the maximal logarithmic amplitude, and excludes the data below that range.



- *mirenvelope(..., 'Mu', mu)* computes the logarithm of the envelope, before the eventual differentiation, using a mu-law compression (Klapuri et al., 2006). Default value for *mu* when 'Mu' is called: 100
- *mirenvelope(..., 'Power')* computes the power (square) of the envelope.
- *mirenvelope(..., 'Diff')* computes the differentiation of the envelope, i.e., the differences between successive samples.



- *mirenvelope(..., 'HalfwaveDiff')* performs a half-wave rectification on the differentiated envelope.



- *mirenvelope(..., 'Normal')* normalizes the values of the envelope by fixing the maximum value to 1.
  - If the audio signal is decomposed into segments, each segment is normalized individually.
  - *mirenvelope(..., 'Normal', 'AcrossSegments')*, on the contrary, normalizes across segments, i.e., with respect to the global maxima across all the segments of that audio signal.
- *mirenvelope(..., 'Lambda', l)* sums the half-wave rectified envelope with the non-differentiated envelope, using the respective weight  $0 < l < 1$  and  $(1-l)$ . (Klapuri et al., 2006).
- *mirenvelope(..., 'Smooth', o)* smooths the envelope using a moving average of order *o*. The default value when the option is toggled on:  $o=30$
- *mirenvelope(..., 'Gauss', o)* smooths the envelope using a gaussian of standard deviation *o* samples. The default value when the option is toggled on:  $o=30$

## PRESELECTED MODEL

Complete (or nearly complete) model is available:

- *mirenvelope*(..., '**Klapurio6**') follows the model proposed in (Klapuri et al., 2006). It corresponds to

$e = \text{mirenvelope}(\dots, \text{'Spectro'}, \text{'UpSample'}, \text{'Mu'}, \text{'HalfwaveDiff'}, \text{'Lambda'}, .8);$

$\text{mirsum}(e, \text{'Adjacent'}, 10)$

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- '**Time**': the temporal positions of samples (same as '*Pos*'),
- '**DownSampling**': the value of the '*PostDecim*' option,
- '**Halfwave**': whether the envelope has been half-wave rectified (1) or not (0),
- '**Diff**': whether the envelope has been differentiated (1) or not (0),
- '**Centered**': whether the envelope is centered (1) or not (0),
- '**Phase**': the phase of the spectrogram, if necessary.

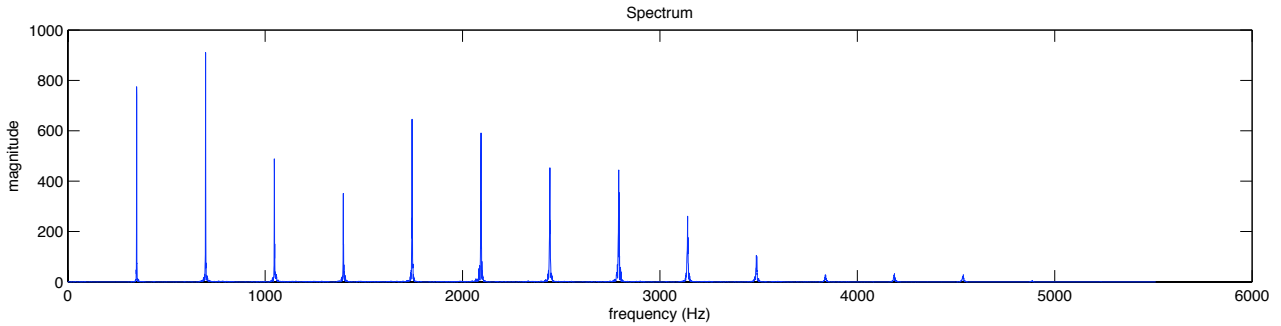
## ***mirspectrum***

### FOURIER TRANSFORM

A decomposition of the signal (be it an audio waveform, or an envelope, etc.) along frequencies can be performed using a Discrete Fourier Transform, which, for an audio signal  $x$  has for equation:

$$X_k = \left| \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N} kn} \right|, k = 0, \dots, N/2$$

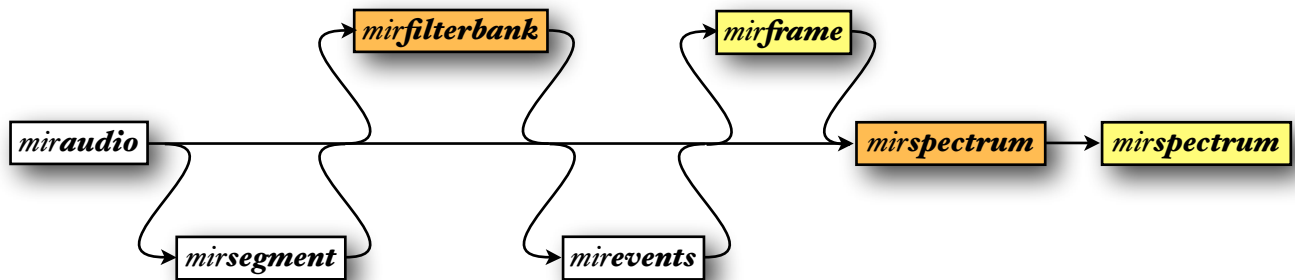
This decomposition is performed using a Fast Fourier Transform by the *mirspectrum* function by calling Matlab *fft* function. The graph returned by the function highlights the repartition of the amplitude of the frequencies (i.e., the modulus of  $Xk$  for all  $k$ ), such as the following:



We can also obtain for each frequency the actual phase position (i.e., the phase of  $Xk$ ), which indicates the exact position of each frequency component at the instant  $t = 0$ . If the result of the spectrum decomposition is  $s$ , the phase spectrum is obtained by using the command:

*get(s, 'Phase')*

### FLOWCHART INTERCONNECTIONS



*mirspectrum* accepts as input data type either:

- ***miraudio*** objects, where the audio waveform can be segmented (using ***mirsegment***), decomposed into channels (using ***mirfilterbank***), and/or decomposed into frames (using ***mirframe***);

- **file name** or the **'Folder'** keyword;
- data in the detection curve category (cf. *mirevents*):
  - *mirenvelope* objects, frame-decomposed or not,
  - fluxes (cf. *mirflux*), frame-decomposed or not;
- *mirspectrum* frame-decomposed objects, with frequency redistributed into bands ('Mel', 'Bark', cf. below): by calling again *mirspectrum* with the **'AlongBands'** option, Fourier transforms are computed this time on each temporal signal related to each separate band.

## FRAME DECOMPOSITION

*mirspectrum*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. 'Frame' section.

## PARAMETERS SPECIFICATION

The range of frequencies, in Hz, can be specified by the options:

- *mirspectrum*(..., **'Min'**, *mi*) indicates the lowest frequency taken into consideration, expressed in Hz. Default value: 0 Hz.
- *mirspectrum*(..., **'Max'**, *ma*) indicates the highest frequency taken into consideration, expressed in Hz. Default value: the maximal possible frequency, corresponding to the sampling rate divided by 2.
- *mirspectrum*(..., **'Window'**, *w*) specifies the windowing method. Windows are used to avoid the problems due to the discontinuities provoked by finite signals. Indeed, an audio sequence is not infinite, and the application of the Fourier Transform requires to replace the infinite time before and after the sequence by zeroes, leading to possible discontinuities at the borders. Windows are used to counteract those discontinuities. Possible values for *w* are either *w* = 0 (no windowing) or any windowing function proposed in the *Signal Processing Toolbox*<sup>1</sup>. Default value: *w* = 'hamming', the Hamming window being a particular good window for Fourier Transform.
- *mirspectrum*(..., **'NormalInput'**) normalizes the waveform between 0 and 1 before computing the Fourier Transform.

---

<sup>1</sup> The list of possible window arguments can be found in the *window* documentation (*help window*).

- *mirspectrum*(..., '**Phase**', 'No') does not compute the related FFT phase. The FFT phase is not computed anyway whenever another option that will make the phase information irrelevant (such as '*Log*', '*dB*', etc.) is specified.

## RESOLUTION SPECIFICATION

The frequency resolution of the spectrum directly depends on the size of the audio waveform: the longer the waveform, the better the frequency resolution. It is possible, however, to (artificially) increase the frequency resolution of a given audio waveform by simply adding a series of zeros at the end of the sequence, which is called *zero-padding*.

- *mirspectrum*(..., '**Power2**', *p*) forces the length of the input signal to be a power of 2, through zero-padding. This was known to enable to speed up the FFT computation. This constraint is toggled on by default.
- *mirspectrum*(..., '**MinRes**', *mr*) adds a constraint related to the a minimal frequency resolution, fixed to the value *mr* (in Hz). If '*Power2*' is also toggled on, the audio waveform is automatically zero-padded to the lowest power of 2 ensuring the required frequency resolution.
- *mirspectrum*(..., '*MinRes*', *mr*, '**OctaveRatio**', *tol*): Indicates the minimal accepted resolution in terms of number of divisions of the octave. Low frequencies are ignored in order to reach the desired resolution. The corresponding required frequency resolution is equal to the difference between the first frequency bins, multiplied by the constraining multiplicative factor *tol* (set by default to .75).
- *mirspectrum*(..., '**Res**', *r*) specifies the frequency resolution *r* (in Hz) that will be secured as closely as possible, through an automated zero-padding. The length of the resulting audio waveform will not necessarily be a power of 2, i.e., the '*Power2*' constraint is not taken into consideration.
- *mirspectrum*(..., '**Length**', *l*) specifies the length of the audio waveform after zero-padding. The length of the resulting audio waveform will not necessarily be a power of 2, i.e., the '*Power2*' constraint is not taken into consideration.
- *mirspectrum*(..., '**ZeroPad**', *s*) performs a zero-padding of *s* samples. The length of the resulting audio waveform will not necessarily be a power of 2, i.e., the '*Power2*' constraint is not taken into consideration.
- *mirspectrum*(..., '**WarningRes**', *mr*) indicates a required frequency resolution, in Hz, for the input signal. If the resolution does not reach that prerequisite, a warning is displayed.

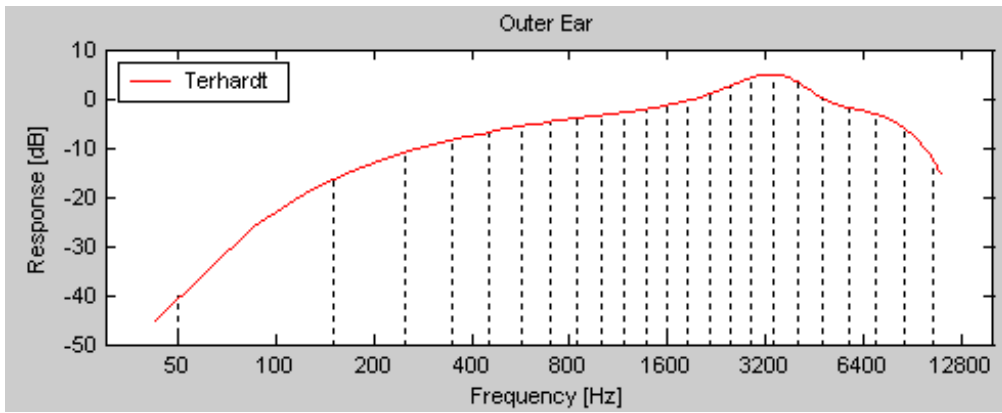
Alternatively, the spectrum decomposition can be performed through a Constant Q Transform instead of a FFT, which enables to express the frequency resolution as a constant number of bins per octave:

- *mirspectrum*(..., '**ConstantQ**', *nb*) fixes the number of bins per octave to *nb*. Default value when the '**ConstantQ**' option is toggled on: *nb*=12 bins per octave.

Please note however that the Constant Q Transform is implemented as a *Matlab* M file, whereas *Matlab*'s FFT algorithm is optimized, therefore faster.

## POST-PROCESSING OPTIONS

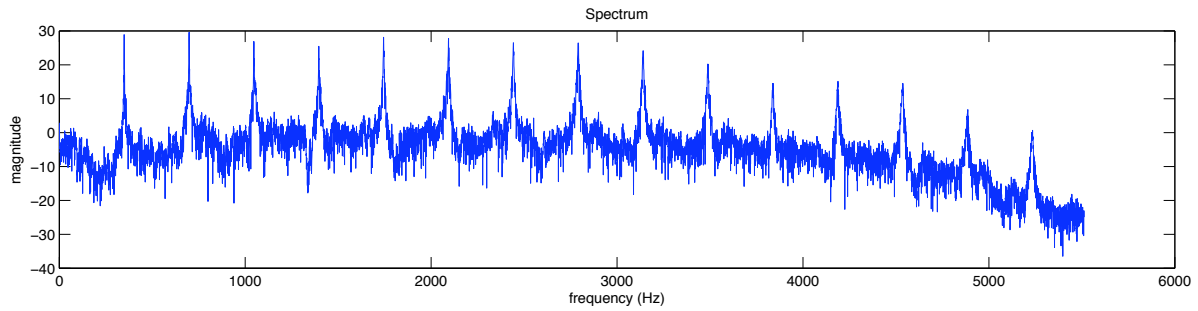
- *mirspectrum*(..., '**Terhardt**') applies Terhardt (1979) outer ear model. The function is mainly characterized by an attenuation in the lower and higher registers of the spectrum, and an emphasis around 2-5 KHz, where much of the speech information is carried. (Code based on Pampalk's MA toolbox).



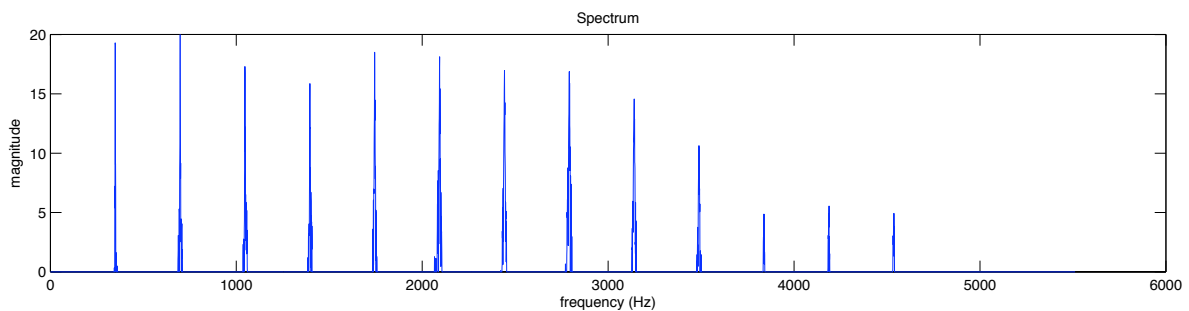
- *mirspectrum*(..., '**NormalLength**') normalises with respect to the duration (in s.) of the audio input data.
- *mirspectrum*(..., '**Power**') computes the *power spectral density* (PSD):

$$X_k = \left| \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N} kn} \right|^2, k = 0, \dots, N/2$$

- *mirspectrum*(..., '**Normal**', *n*) normalises with respect to the *power spectrum*, i.e., the summation of the PSD over all frequencies. Toggled off by default. Two normalisation strategies are offered:
  - *n* = '**Local**' normalises each frame separately. Default value when the '**Normal**' option is toggled on.
  - *n* = '**Global**' normalises with respect to the global maximum over all frames and channels.
- *mirspectrum*(..., '**dB**') represents the power spectrum in decibel scale. For the previous example we obtain the following spectrum:



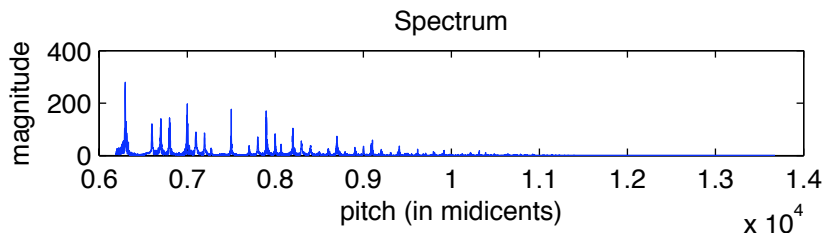
- `mirspectrum(..., 'dB', th)` keeps only the highest energy over a range of *th* dB. The highest energy is computed across all frames. For example if we take only the 20 most highest dB in the previous example we obtain:



- `mirspectrum(..., 'Resonance', r)` multiplies the spectrum curve with a resonance curve that emphasizes pulsations that are more easily perceived. Two resonance curves are available:
  - $r = \text{'ToivaiainenSnyder'}$  (Toivaiainen & Snyder 2003), default choice, used for event detection (cf. `mirtempo`),
  - $r = \text{'Fluctuation'}$ : fluctuation strength (Fastl 1982), default choice for frame-decomposed `mirspectrum` objects redecomposed in 'Mel' bands (cf. `mirfluctuation`).
- `mirspectrum(..., 'Smooth', o)` smooths the spectrum curve using a moving average of order *o*. Default value when the option is toggled on:  $o=10$
- `mirspectrum(..., 'Gauss', o)` smooths the spectrum curve using a gaussian of standard deviation *o* samples. Default value when the option is toggled on:  $o=10$
- `mirspectrum(..., 'TimeSmooth', o)` smooths each frequency channel of a spectrogram using a moving average of order *o*. Default value when the option is toggled on:  $o=10$
- `mirspectrum(..., 'MinGate', th)` turns all value below the threshold into zero values. For linear spectrum, the threshold is expressed as a ratio between 0 and 1 with respect to the global maximum. For logarithm spectrum, the threshold is expressed as a distance with respect to the global maximum, and the zero values are actually -Inf values. For instance, for a spectrum in dB,  $th = 20$  indicates that value under 20 db below the global maximum are turned into -Inf values.

## FREQUENCY REDISTRIBUTION

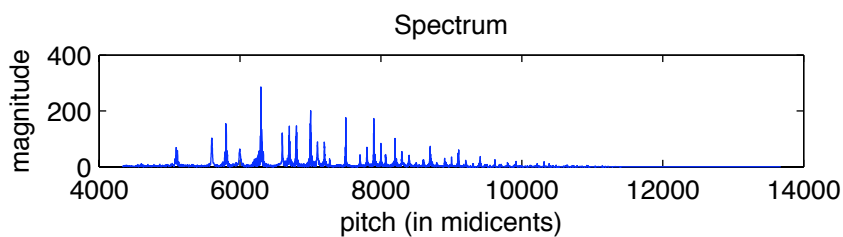
- *mirspectrum(..., 'Cents')* redistributes the frequencies along cents. Each octave is decomposed into 1200 bins equally distant in the logarithmic representation. The frequency axis is hence expressed in MIDI-cents unit: to each pitch of the equal temperament is associated the corresponding MIDI pitch standard value multiply by 100 ( $69 \times 100 = 6900$  for A4=440Hz,  $70 \times 100 = 7000$  for B4, etc.).



*mirspectrum('ragtime.wav','Cents')*

It has to be noticed that this decomposition requires a frequency resolution that gets higher for lower frequencies: a cent-distribution starting from infinitely low frequency (near 0 Hz would require an infinite frequency resolution). Hence by default, the cent-decomposition is defined only for the frequency range suitable for the frequency resolution initially associated to the given spectrum representation. Two levers are available here:

- If a minimal frequency range for the spectrum representation has been set (using the *'Min'* parameter), the frequency resolution of the spectrum is automatically set in order to meet that particular requirement.

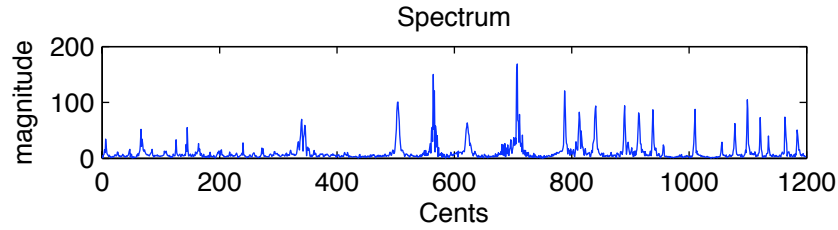


*mirspectrum('ragtime.wav','Cents','Min',100)*

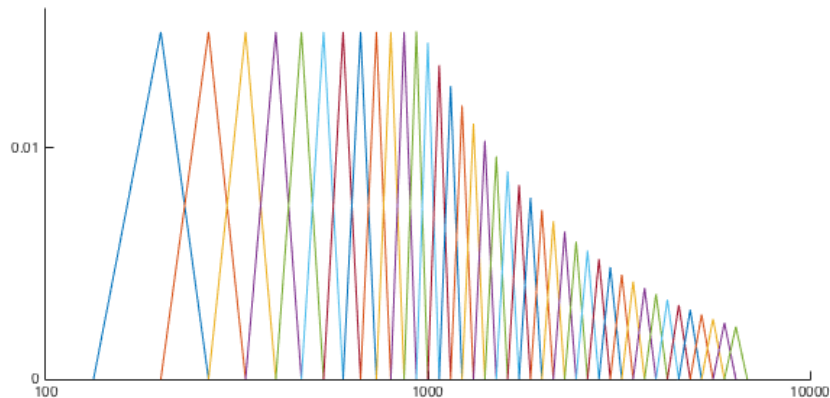
- By increasing the frequency resolution of the spectrum (for instance by using the *'Res'* or *'MinRes'* parameters), the frequency range will be increased accordingly.
- *mirspectrum(..., 'Collapsed')* collapses the cent-spectrum into one octave. In the resulting spectrum, the abscissa contains in total 1200 bins, representing the 1200 cents of one octave, and each bin contains the energy related to one position of one octave and of all the multiple of this octave.

*mirspectrum('ragtime.wav','Cents','Min',100,'Collapsed')*

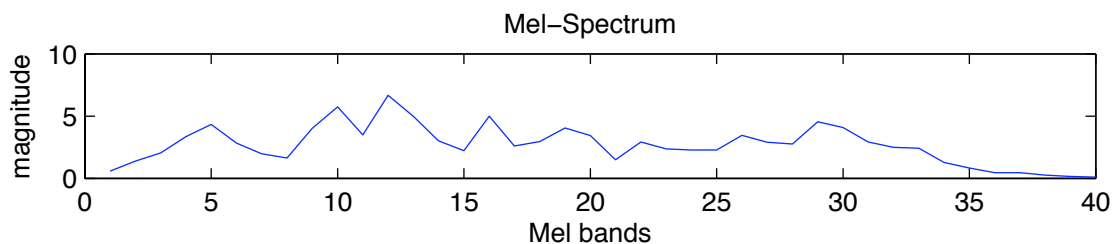




- *mirspectrum(..., 'Mel')* redistributes the frequencies along Mel bands. The Mel-scale of auditory pitch was established on the basis of listening experiments with simple tones (Stevens and Volkman, 1940). The Mel scale is now mainly used for the reason of its historical priority only. It is closely related to the Bark scale. It uses the Mel-scale definition from the *Auditory Toolbox*: 13 linearly-spaced filters (66.66 Hz<sup>2</sup> between center frequencies, with the first one centered around 133.33 Hz) followed by log-spaced filters (separated by a factor of 1.0711703 in frequency.)
- *mirspectrum(..., 'Bands', b)* specifies the number of band in the decomposition. By default *b* = 40.



In our example we obtain the following:



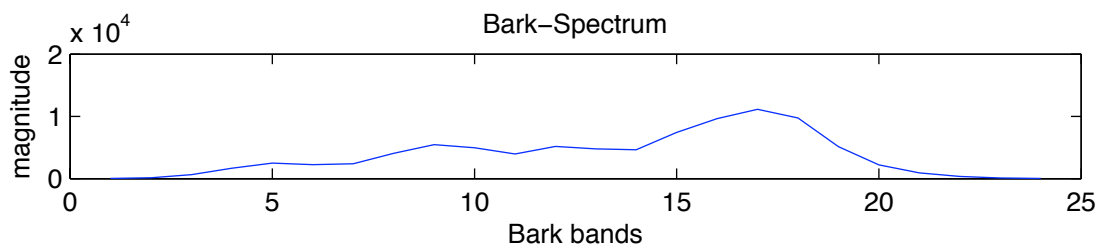
The Mel-scale transformation requires a sufficient frequency resolution of the spectrum: as the lower bands are separated with a distance of 66.66 Hz, the frequency resolution should be hi-

---

<sup>2</sup> The *Auditory Toolbox* documentation erroneously indicates the double: 133.33 Hz between center frequencies. But the figure in the documentation and the actual code uses 66.66 Hz.

gher than 66.66 Hz in order to ensure that each Mel band can be associated with at least one frequency bin of the spectrum. If the ‘*Mel*’ option is performed in the same *mirspectrum* command that performs the actual FFT, then the minimal frequency resolution is implicitly ensured, by forcing the minimal frequency resolution (‘*MinRes*’ parameter) to be equal or below 66 Hz. If on the contrary the ‘*Mel*’ is performed in a second step, and if the frequency resolution is worse than 66 Hz, then a warning message is displayed in the Command Window.

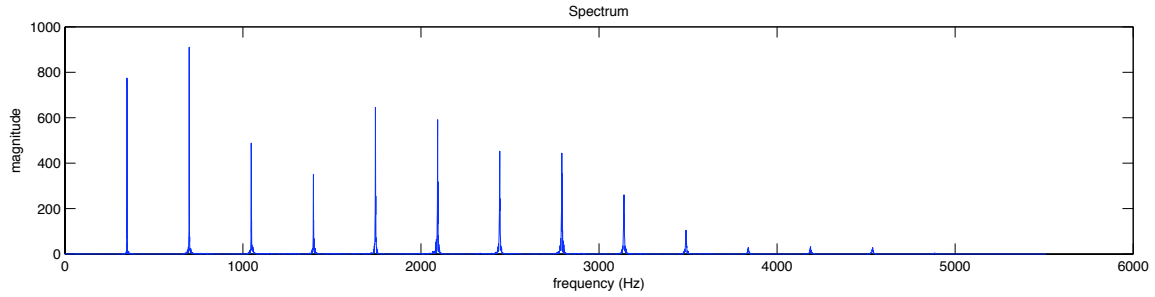
- *mirspectrum*(..., ‘**Bark**’) redistributes the frequencies along critical band rates (in Bark). Measurement of the classical "critical bandwidth" typically involves loudness summation experiments (Zwicker et al., 1957). The critical band rate scale differs from Mel-scale mainly in that it uses the critical band as a natural scale unit. The code is based on the *MA* toolbox.
- *mirspectrum*(..., ‘**Mask**’) models masking phenomena in each band: when a certain energy appears at a given frequency, lower frequencies in the same frequency region may be unheard, following particular equations. By modeling these masking effects, the unheard periodicities are removed from the spectrum. The code is based on the *MA* toolbox. In our example this will lead to:



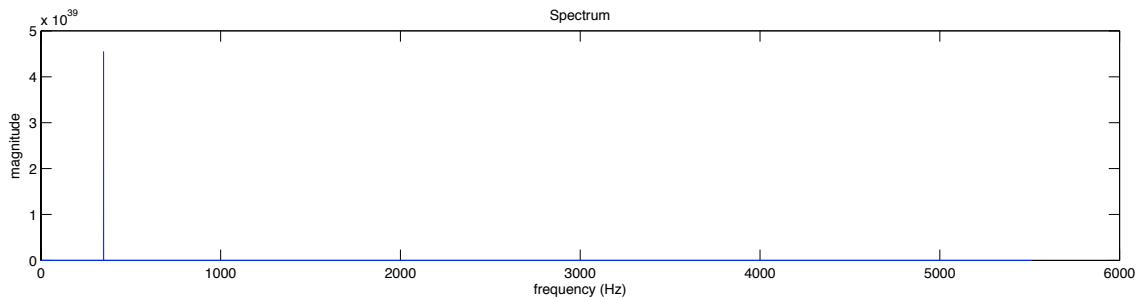
## HARMONIC SPECTRAL ANALYSIS

A lot of natural sounds, especially musical ones, are harmonic: each sound consists of a series of frequencies at a multiple ratio of the one of lowest frequency, called fundamental. Techniques have been developed in signal processing to reduce each harmonic series to its fundamental, in order to simplify the representation. *MIRtoolbox* includes two related techniques for the attenuation of harmonics in spectral representation (Alonso et al, 2003):

- *mirspectrum*(..., ‘**Prod**’, *m*) Enhances components that have harmonics located at multiples of range(s) *m* of the signal's fundamental frequency. Computed by compressing the signal by a list of factors *m*, and by multiplying all the results with the original signal. Default value is *m* = 1:6. Hence for this initial spectrum:



we obtain this reduced spectrum:



- *mirspectrum*(..., '**Sum**', *m*) Similar idea using addition of the multiples instead of multiplication.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

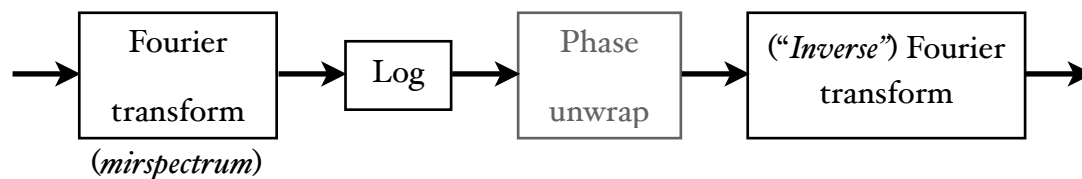
- '**Frequency**': the frequency (in Hz.) associated to each bin (same as '*Pos*'),
- '**Magnitude**': the magnitude associated to each bin (same as '*Data*'),
- '**Phase**': the phase associated to each bin,
- '**XScale**': whether the frequency scale has been redistributed into cents – with ('*Cents(Collapsed)*') or without ('*Cents*') collapsing into one octave –, mels ('*Mel*'), barks ('*Bark*'), or not redistributed at all ('*Freq*'),
- '**Power**': whether the spectrum has been squared (i) or not (o),
- '**Log**': whether the spectrum is in log-scale (i) or in linear scale (o).

## *mircepstrum*

### SPECTRAL ANALYSIS OF SPECTRUM

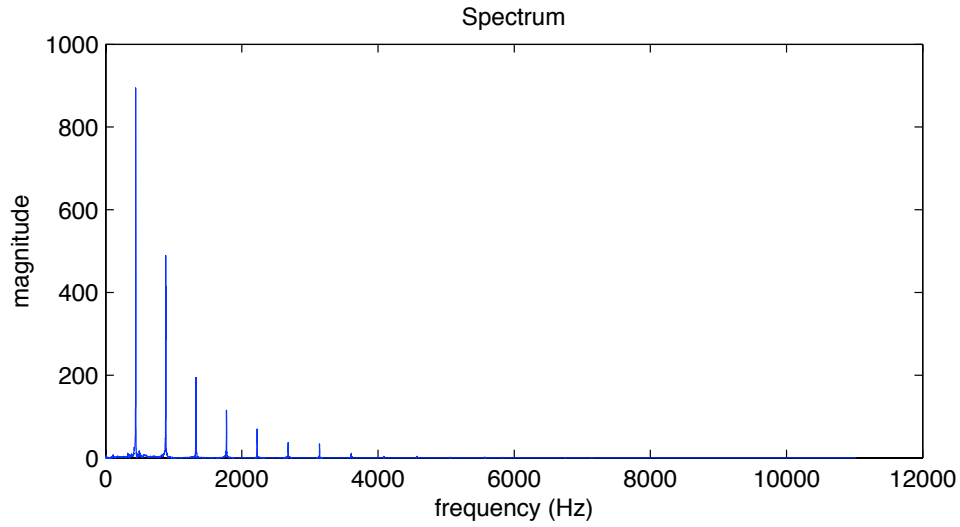
The harmonic sequence can also be used for the detection of the fundamental frequency itself. One idea is to look at the spectrum representation, and try to automatically detect these periodic sequences. And one simple idea consists in performing a Fourier Transform of the Fourier Transform itself, leading to a so-called cepstrum (Bogert et al., 1963).

So if we take the complex spectrum ( $Xk$  in the equation defining *mirpectrum*), we can operate the following chain of operations:

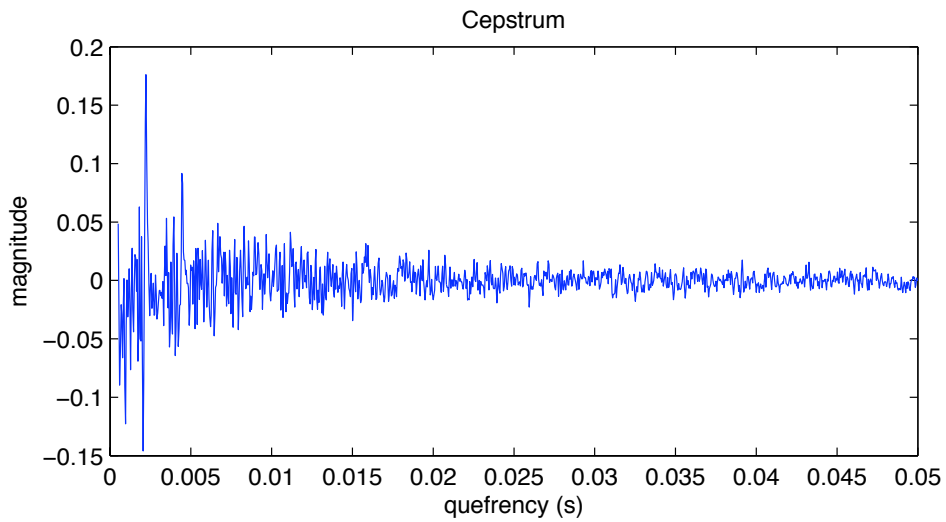


- First a logarithm is performed in order to allow an additive separability of product components of the original spectrum. For instance, for the voice in particular, the spectrum is composed of a product of a vocal cord elementary burst, their echoes, and the vocal track. In the logarithm representations, these components are now added one to each other, and we will then be able to detect the periodic signal as one of the components.
- Then because the logarithm provokes some modification of the phase, it is important to ensure that the phase remains continuous.
- Finally the second Fourier transform is performed in order to find the periodic sequences. As it is sometime a little difficult to conceive what a Fourier transform of Fourier transform is really about, we can simply say, as most say, that it is in fact an Inverse Fourier Transform (as it is the same thing, after all), and the results can then be expressed in a kind of temporal domain, with unit called "quefreny".

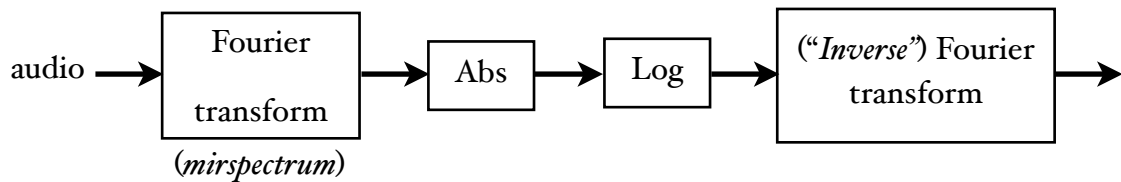
For instance for this spectrum:



we obtain the following cepstrum:

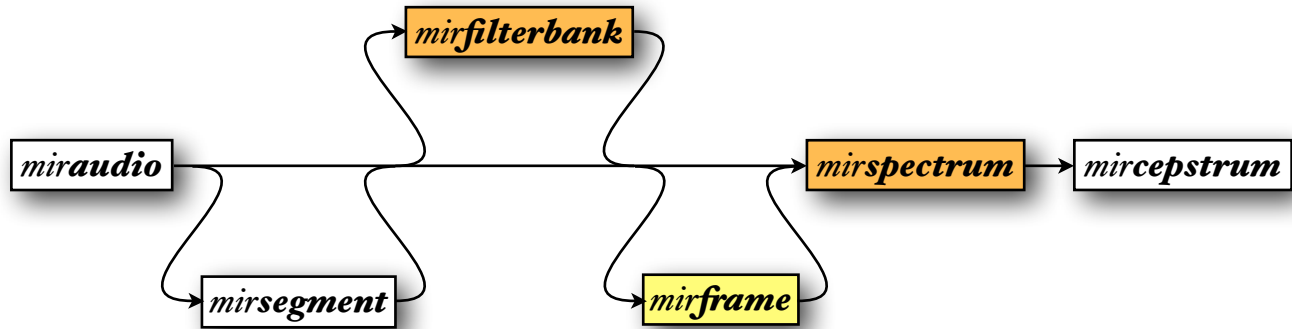


The cepstrum can also be computed from the spectrum amplitude only, by simply taking the logarithm, and directly computing the Fourier transform.



In this case, the phase of the spectrum is not computed.

## FLOWCHART INTERCONNECTIONS



*mircepstrum* accepts either:

- *mirspectrum* objects, or
- *miraudio* objects (same as for *mirspectrum*),
- **file name** or the **'Folder'** keyword.

## FRAME DECOMPOSITION

*mircepstrum*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. **'Frame'** section.

## PARAMETER SPECIFICATIONS

- *mircepstrum*(..., **'Freq'**): The results can be represented, instead of using the quefrency domain (in seconds), back to the frequency domain (in Hz) by taking the inverse of each abscissae value. In this frequency representation, each peak is located on a position that directly indicates the associated fundamental frequency.
- *mircepstrum*(..., **'Min'**, *min*) specifies the lowest delay taken into consideration, in seconds. Default value: 0.0002 s (corresponding to a maximum frequency of 5 kHz). This default value is not set to 0 s in order to exclude the very high peaks confined in the lowest quefrency region: these high peaks seem to come from the fact that the spectrum is a non-centered signal, thus with high (quasi-)stationary energy. However, the value can be forced to 0 using this **'Min'** option.
- *mircepstrum*(..., **'Max'**, *max*) specifies the highest delay taken into consideration, in seconds. Default value: 0.05 s (corresponding to a minimum frequency of 20 Hz). This default value is not set to *Inf* in order to exclude the very high peaks confined in the highest quefrency region: these high peaks seem to come from the fact that the spectrum is a highly variable signal, thus with high energy on its highest frequencies. However, the value can be forced to *Inf* using this **'Max'** option.

- *mircepstrum*(..., '**Complex**') computes the cepstrum using the complex spectrum. By default, the cepstrum is computed from the spectrum amplitude only.

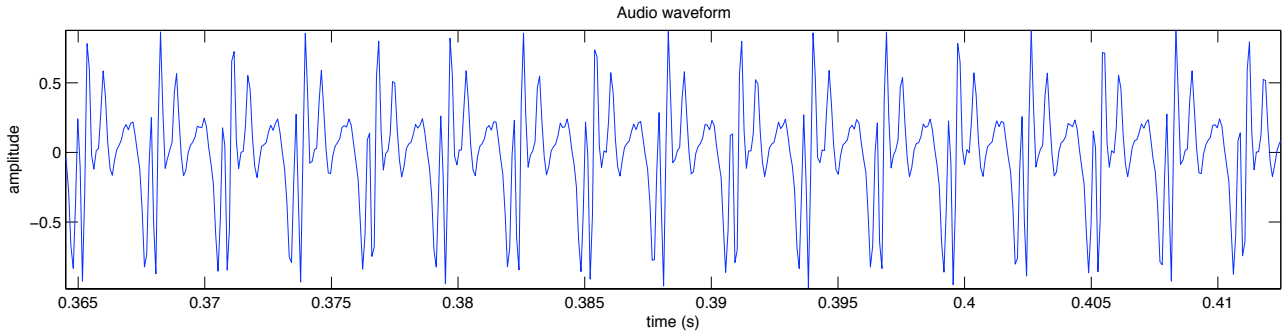
## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- '**Magnitude**': same as '*Data*'
- '**Phase**': the phase related to the magnitude (only if '*Complex*' option is used)
- '**FreqDomain**': whether the axis data are frequencies in s. (o) or frequency in Hz. (i)
- To obtain the frequency or quefrency values associated to each bin, we recommend using '*Pos*', in conjunction with '*FreqDomain*'.

## AUTOCORRELATION FUNCTION

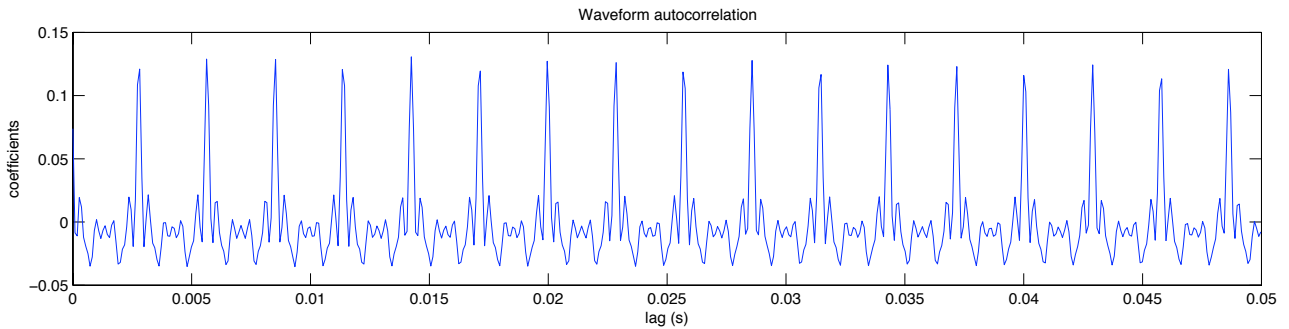
Another way to evaluate periodicities in signals (be it an audio waveform, a spectrum, an envelope, etc.) consists in looking at local correlation between samples. If we take a signal  $x$ , such as for instance this trumpet sound:



the autocorrelation function is computed as follows:

$$R_{xx}(j) = \sum_n x_n \bar{x}_{n-j} .$$

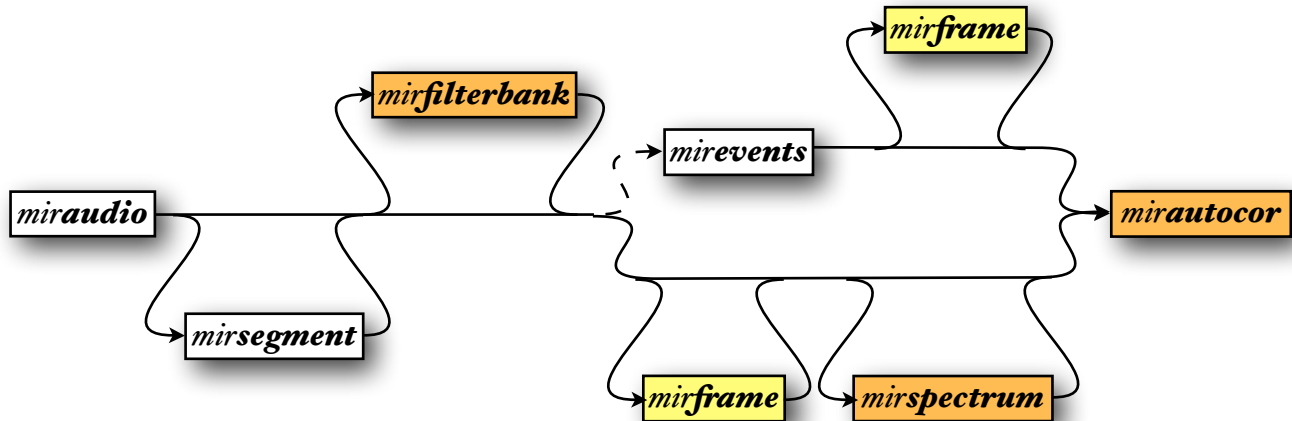
For a given lag  $j$ , the autocorrelation  $R_{xx}(j)$  is computed by multiplying point par point the signal with a shifted version of it of  $j$  samples. We obtain this curve:



Hence when the lag  $j$  corresponds to a period of the signal, the signal is shifted to one period ahead, and therefore is exactly superposed to the original signal. Hence the summation gives very high value, as the two signals are highly correlated.



## FLOWCHART INTERCONNECTIONS



*mirautocor* usually accepts either:

- **file name** or the **'Folder'** keyword,
- ***miraudio*** objects, where the audio waveform can be segmented (using ***mirsegment***), decomposed into channels (using ***mirfilterbank***), and/or decomposed into frames (using ***mirframe***),
- ***mirspectrum*** objects,
- data in the event detection curve category (cf. ***mirevents***):
  - ***mirenvelope*** objects, frame-decomposed or not,
  - fluxes (cf. ***mirflux***), frame-decomposed or not,
- ***mirautocor*** objects, for further processing.

## FRAME DECOMPOSITION

*mirautocor*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. **'Frame'** section.

## PARAMETERS SPECIFICATION

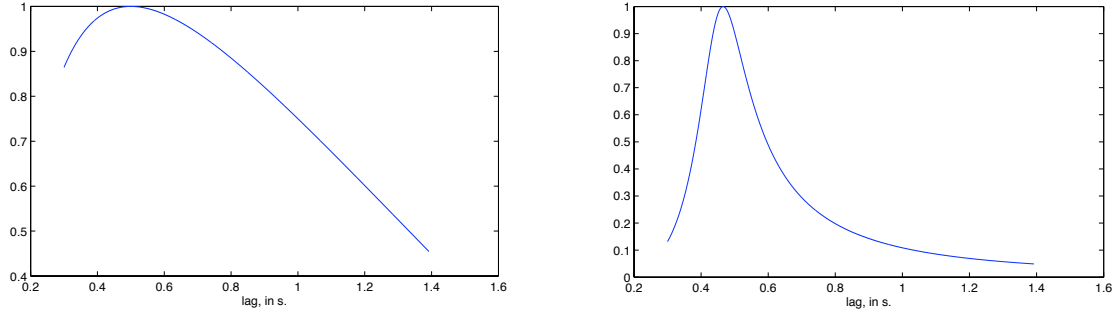
- *mirautocor*(..., **'Min'**, *mi*) indicates the lowest delay taken into consideration. Default value: 0 s. The unit can be specified:
  - *mirautocor*(..., **'Min'**, *mi*, **'s'**) (default unit)
  - *mirautocor*(..., **'Min'**, *mi*, **'Hz'**)

- *mirautocor*(..., '**Max**', *ma*) indicates the highest delay taken into consideration. The unit can be specified as for '*Min*'. Default value:
  - if the input is an audio waveform, the highest delay is by default 0.05 s (corresponding to a minimum frequency of 20 Hz).
  - if the input is an envelope, the highest delay is by default 2 s.
- *mirautocor*(..., '**Normal**', *n*) specifies a normalization option for the cross-correlation ('*biased*', '*unbiased*', '*coeff*', '*none*'). This corresponds exactly to the normalization options in *Matlab xcorr* function, as *mirautocor* actually calls *xcorr* for the actual computation. The default value is '*coeff*', corresponding to a normalization so that the autocorrelation at zero lag is identically 1. If the data is multi-channel, the normalization is such that the sum over channels at zero lag becomes identically 1. Note however that the '*coeff*' routine is not used when the compression ('*Compres*') factor *k* is not equal to 2 (see below).

## POST-PROCESSING OPTIONS

- *mirautocor*(..., '**Freq**') represents the autocorrelation function in the frequency domain: the periods are expressed in Hz instead of seconds (see the last curve in the figure below for an illustration).
- *mirautocor*(..., '**NormalWindow**') divides the autocorrelation by the autocorrelation of the window. Boersma (1993) shows that by default the autocorrelation function gives higher coefficients for small lags, since the summation is done on more samples. Thus by dividing by the autocorrelation of the window, we normalize all coefficients in such a way that this default is completely resolved. At first sight, the window should simply be a simple rectangular window. But Boersma (1993) shows that it is better to use '*banning*' window in particular, in order to obtain better harmonic to noise ratio.
  - *mirautocor*(..., '**NormalWindow**', *w*) specifies the window to be used, which can be any window available in the *Signal Processing Toolbox*. Besides *w* = '*rectangular*' will not perform any particular windowing (corresponding to a rectangular ("invisible") window), but the normalization of the autocorrelation by the autocorrelation of the invisible window will be performed nonetheless. The default value is *w* = '*banning*'.
  - *mirautocor*(..., '**NormalWindow**', '*off*') toggles off this normalization (which is '*on*' by default).
- *mirautocor*(..., '**Resonance**', *r*) multiplies the autocorrelation curve with a resonance curve that emphasizes pulsations that are more easily perceived. Two resonance curves are proposed:
  - *r* = '*ToiviainenSnyder*' (Toiviainen & Snyder 2003) (default value if '*Resonance*' option toggled on),

- $r = \text{'vanNoorden'}$  (van Noorden & Moelants, 1999).



Resonance curves 'ToiviainenSnyder' (left) and 'vanNoorden' (right)

This option should be used only when the input of the *mirautocor* function is an amplitude envelope, i.e., a *mirenvelope* object.

- *mirautocor*(..., **'Center'**,  $c$ ) assigns the center value of the resonance curve, in seconds. Works mainly with 'ToiviainenSnyder' option. Default value:  $c = 0.5$ .
- *mirautocor*(..., **'Halfwave'**) performs a half-wave rectification on the result, in order to just show the positive autocorrelation coefficients.

## GENERALIZED AUTOCORRELATION

*mirautocor*(..., **'Compres'**,  $k$ ) – or equivalently *mirautocor*(..., **'Generalized'**,  $k$ ) – computes the autocorrelation in the frequency domain and includes a magnitude compression of the spectral representation. Indeed an autocorrelation can be expressed using Discrete Fourier Transform as

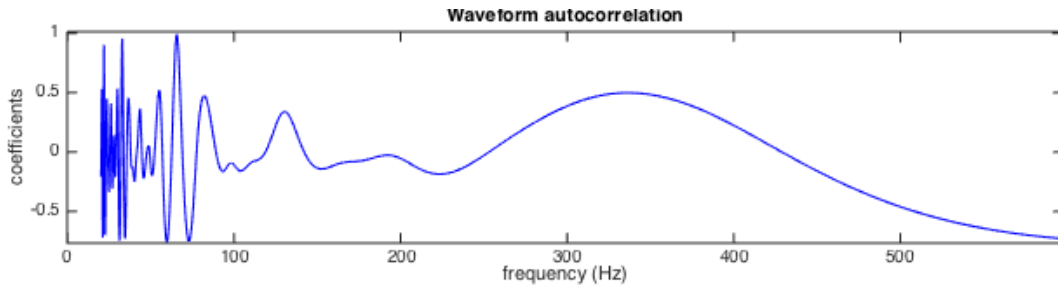
$$y = IDFT(|DFT(x)|^2),$$

which can be generalized as:

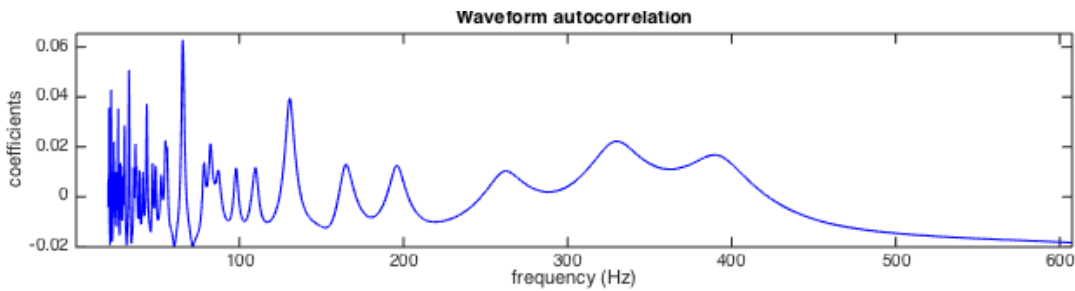
$$y = IDFT(|DFT(x)|^k),$$

Compression of the autocorrelation (i.e., setting a value of  $k$  lower than 2) are recommended in (Tolonen & Karjalainen, 2000) because this decreases the width of the peaks in the autocorrelation curve, at the risk however of increasing the sensitivity to noise. According to this study, a good compromise seems to be achieved using value  $k = .67$ . By default, no compression is performed (hence  $k = 2$ ), whereas if the 'Compres' keyword is used, value  $k = .67$  is set by default if

no other value is indicated.



*mirautocor('Cmaj.wav', 'Freq')*



*mirautocor('Cmaj.wav', 'Freq', 'Compress')*

## ENHANCED AUTOCORRELATION

In the autocorrelation function, for each periodicity in the signal, peaks will be shown not only at the lag corresponding to that periodicity, but also to all the multiples of that periodicity. In order to avoid such redundancy of information, techniques have been proposed that automatically remove these harmonics. In the frequency domain, this corresponds to sub-harmonics of the peaks.

*mirautocor(..., 'Enhanced', a)*: The original autocorrelation function is half-wave rectified, time-scaled by factor  $a$  (which can be a factor list as well), and subtracted from the original clipped function (Tolonen & Karjalainen, 2000). If the 'Enhanced' option is not followed by any value, the default value is  $a = 2:10$ , i.e., from 2 to 10.

If the curve does not start from zero at low lags but begins instead with strictly positive values, the initial discontinuity would be propagated throughout all the scaled version of the curve. In order to avoid this phenomenon, the curve is modified in two successive ways:

- if the curve starts with a descending slope, the whole part before the first local minimum is removed from the curve,
- if the curve starts with an ascending slope, the curve is prolonged to the left following the same slope but which is increased by a factor of 1.1 at each successive bin, until the curve reaches the  $x$ -axis.

See the figure below for an example of enhanced autocorrelation when computing the pitch content of a piano *Amin3* chord, with the successive step of the default enhancement, as used by default in *mirpitch* (cf. description of *mirpitch*).

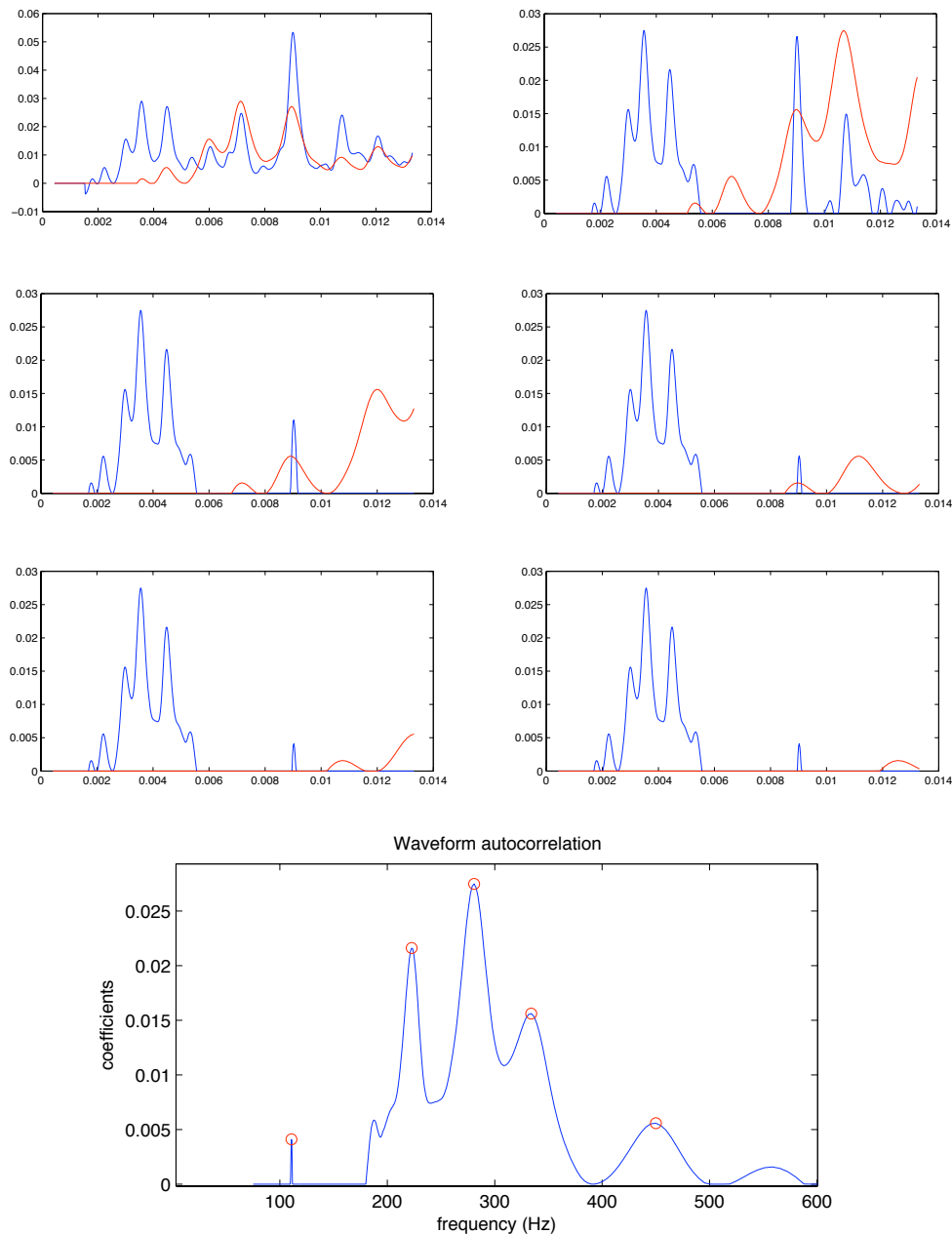
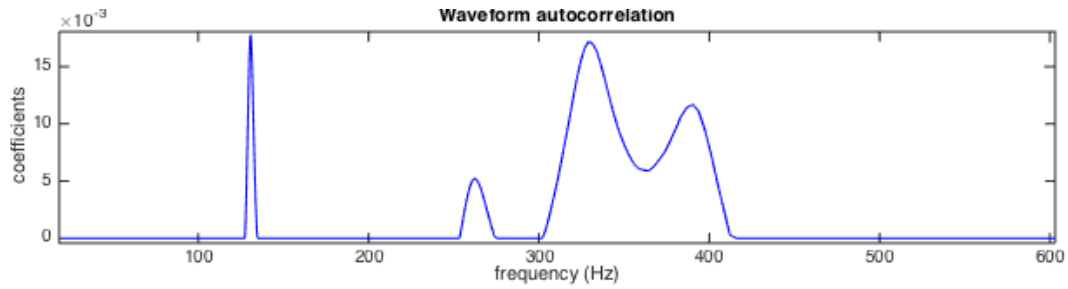


fig 1: Waveform autocorrelation of a piano chord *Amaj3* (blue), and scaled autocorrelation of factor 2 (red);  
 fig 2: subtraction of the autocorrelation by the previous scaled autocorrelation (blue), scaled autocorrelation  
 of factor 3 (red); fig 3: resulting subtraction (blue), scaled autocorrelation of factor 4 (red); fig 4: idem for  
 factor 5; fig 5: idem for factor 6; fig 6: idem for factor 7; fig 7: resulting autocorrelation curve in the frequen-  
 cy domain and peak picking



*mirautocor('Cmaj.wav', 'Freq', 'Compress', 'Enhanced', 2:20)*

## ACCESSIBLE OUTPUT

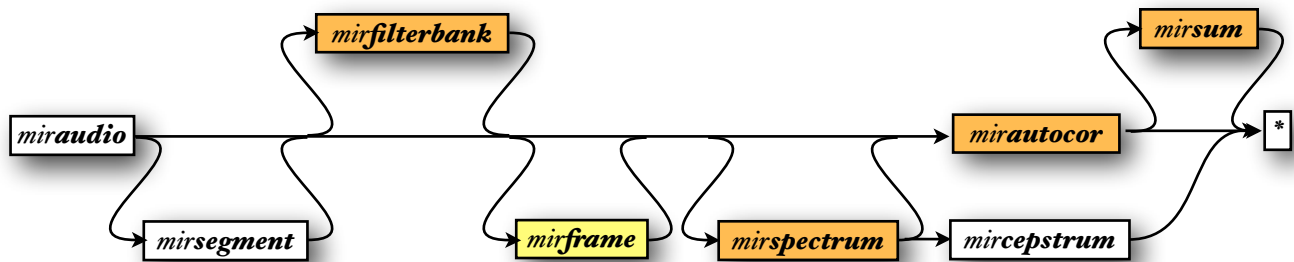
cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **'Coeff'**: the autocorrelation coefficients (same as *'Data'*),
- **'Lag'**: the lags (same as *'Pos'*),
- **'FreqDomain'**: whether the lags are in s. (o) or in Hz. (i),
- **'OfSpectrum'**: whether the input is a temporal signal (o), or a spectrum (i),
- **'Window'**: contains the complete envelope signal used for the windowing,
- **'Resonance'**: indicates the resonance curve that has been applied, or empty string if no resonance curve has been applied.

\*

## COMBINING REPRESENTATIONS

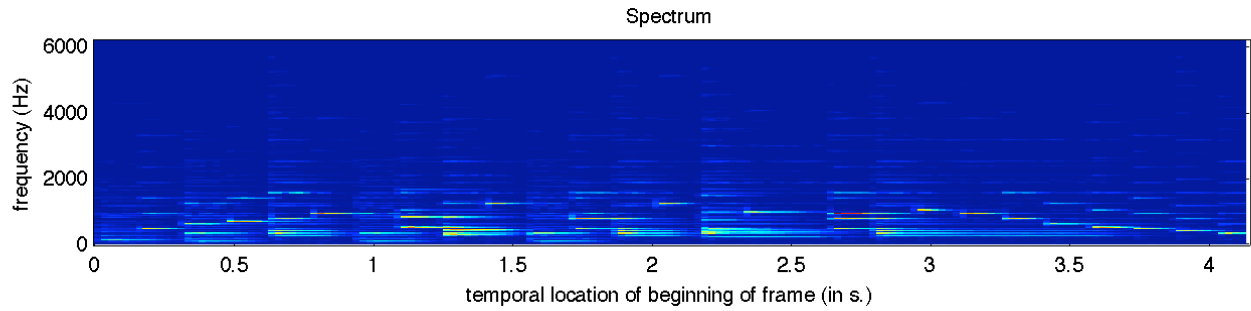
It is also possible to multiple points by points diverse spectral representations and autocorrelation functions, the latter being automatically translated to the spectrum domain (Peeters, 2006). Curves are half-wave rectified before multiplication.



## DISTANCE BETWEEN SUCCESSIVE FRAMES

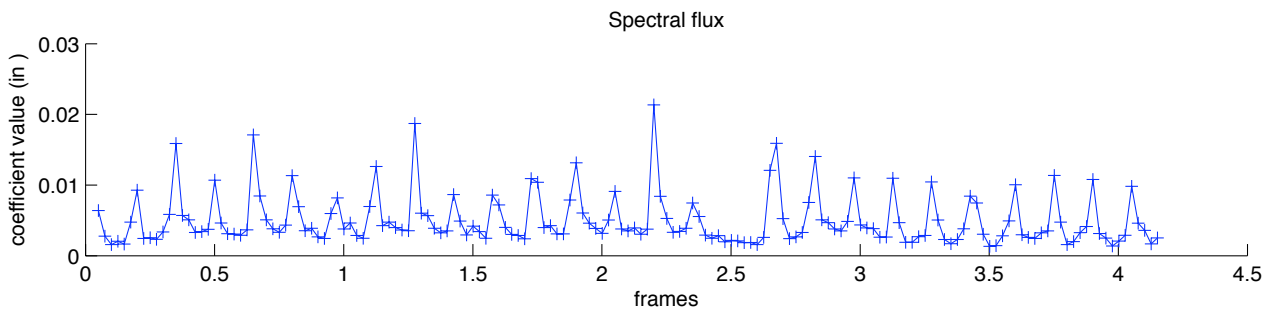
Given a spectrogram:

$$s = \text{mirspectrum}(a, \text{'Frame'})$$



we can compute the spectral flux as being the distance between the spectrum of each successive frames.

$$\text{mirflux}(s)$$

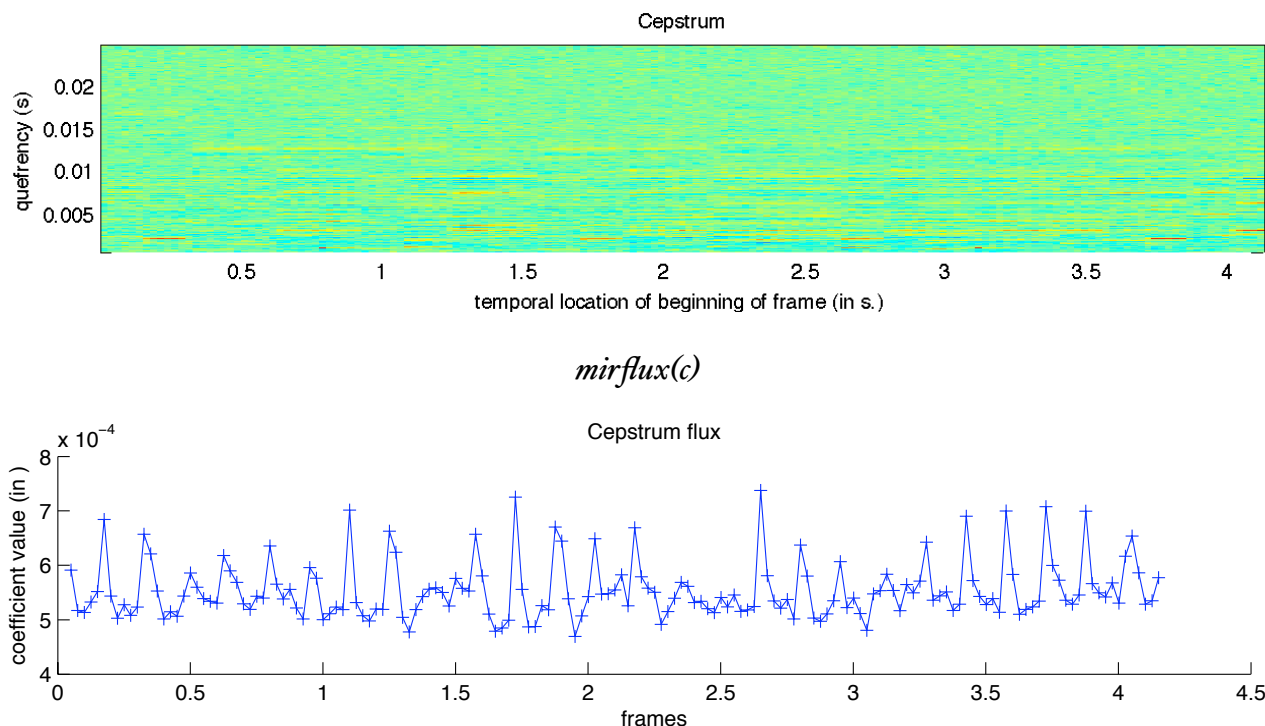


The peaks in the curve indicate the temporal position of important contrast in the spectrogram.

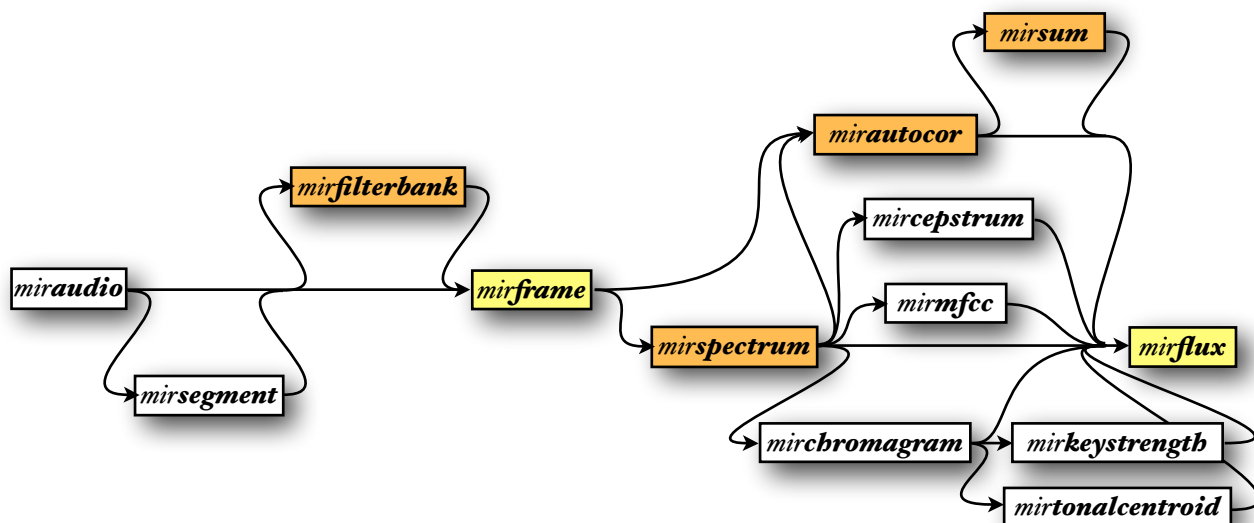
In *MIRtoolbox* fluxes are generalized to any kind of frame-decomposed representation, for instance a cepstral flux:

$$c = \text{mircepstrum}(a, \text{'Frame'})$$





## FLOWCHART INTERCONNECTIONS



*mirflux* usually accepts either:

- ***mirspectrum*** frame-decomposed objects.
- ***miraudio*** objects, where the audio waveform can be segmented (using ***mirsegment***), decomposed into channels (using ***mirfilterbank***). The audio waveform is decomposed into frames if it was not decomposed yet. If the input is a *miraudio* object, the default flux is a spectral flux: i.e., the audio waveform is passed to the *mirspectrum* operator before being fed into *mirflux*.

- **file name** or the **'Folder'** keyword: same behavior than for *miraudio* objects;
- **mirautocor** frame-decomposed objects;
- **mircepstrum** frame-decomposed objects;
- **mirmfcc** frame-decomposed objects;
- **mirchromagram** frame-decomposed objects;
- **mirkeystrength** frame-decomposed objects.

## FRAME DECOMPOSITION

*mirflux*(..., **'Frame'**, ...) specifies the frame configuration, the default being a frame length of 50 ms and half overlapping. For the syntax, cf. the previous *mirframe* vs. *'Frame'* section.

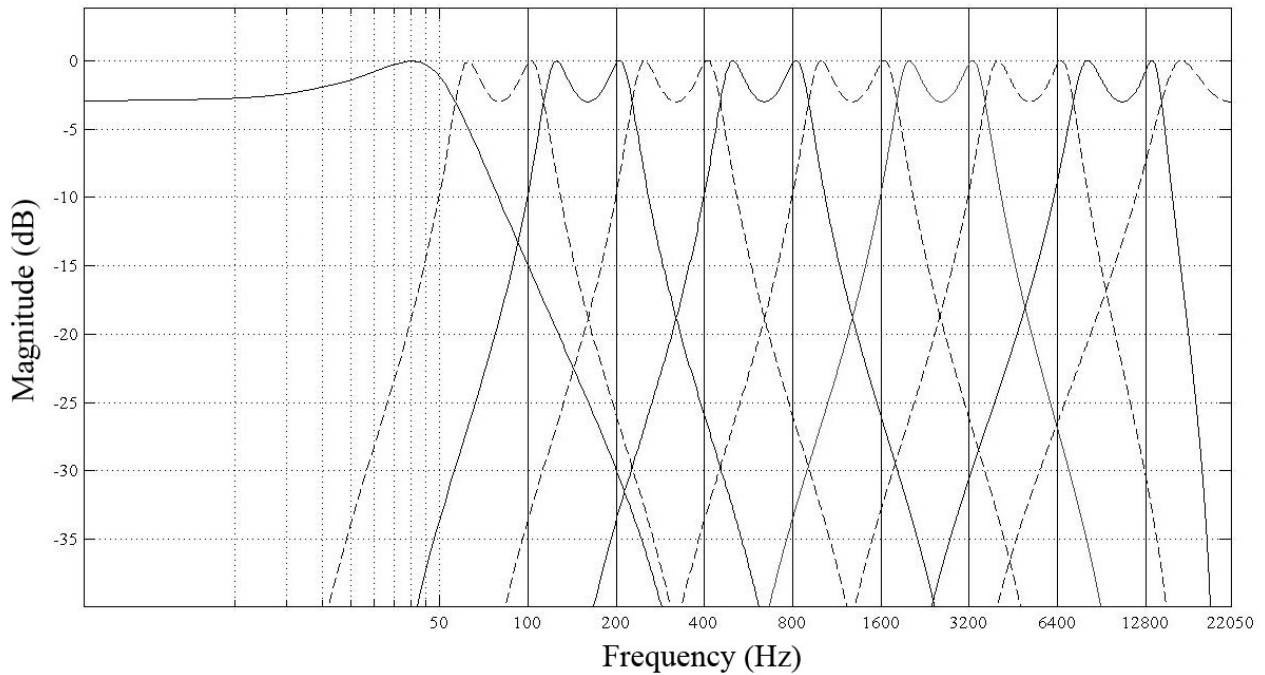
## PARAMETERS SPECIFICATION

- The options **'Min'**, **'Max'**, **'MinGate'** and **'Normal'** from *mirspectrum* can be directly controlled when calling *mirflux*, using the same keywords except for **'Normal'**, which should be replaced by **'NormalSpectrum'** (since there is already a **'Normal'** option in *mirflux*).
- *mirflux*(x, **'Dist'**, *d*) specifies the distance between successive frames, among the list of distances available in *pdist* (cf. *help pdist*). Default distance: **'Euclidean'**.
- *mirflux*(..., **'Inc'**): Only positive difference between frames are summed, in order to focus on increase of energy solely. If toggled on, **'Dist'** parameter can only accept **'Euclidean'** or **'City'**.
- *mirflux*(..., **'Complex'**), for spectral flux, combines the use of both energy and phase information (Bello et al, 2004).
- *mirflux*(..., **'SubBand'**) decomposes first the input waveform using a 10-channel filterbank of octave-scaled second-order elliptical filters, with frequency cut of the first (low-pass) filter at 50 Hz:

*mirfilterbank*(..., **'Manual'**, [-Inf 50\*2.<sup>^(0:1:8)</sup> Inf], **'Order'**, 2)

Alluri and Toivainen (2010, 2012) introduce this model, using a frame size of 25 ms and half-overlapping, corresponding to the command:

*mirflux*(..., **'SubBand'**, **'Frame'**, .025)



The authors found that fluctuations in the lower channels, around 50 Hz - 200 Hz (sub-bands 2 and 3) represent perceived “Fullness”, and those around 1600 - 6400 Hz (sub-Bands 7 and 8) represent perceived “Activity”.

Alternatively, other filterbanks proposed in *mirfilterbank* can be specified, using the syntax:

- *mirflux(..., 'SubBand', 'Gammatone')*
- *mirflux(..., 'SubBand', '2Channels')*

## POST-PROCESSING

- *mirflux(..., 'Median',  $l$ ,  $C$ )*: removes small spurious peaks by subtracting to the result its median filtering. The median filter computes the point-wise median inside a window of length  $l$  (in seconds), that includes a same number of previous and next samples.  $C$  is a scaling factor whose purpose is to artificially rise the curve slightly above the steady state of the signal. If no parameters are given, the default values are:  $l = 0.2$  s. and  $C = 1.3$
- *mirflux(..., 'Median',  $l$ ,  $C$ , 'Halfwave')*: The scaled median filtering is designed to be succeeded by the half-wave rectification process in order to select peaks above the dynamic threshold calculated with the help of the median filter. The resulting signal is called “*detection function*” (Alonso et al., 2003). To ensure accurate detection, the length  $l$  of the median filter must be longer than the average width of the peaks of the detection function.
- *mirflux(..., 'Normal',  $th$ )*: normalises flux. For each flux value (distance between two successive frames  $F_i$  and  $F_{i+1}$ ), the distance is divided by  $\Sigma F_{i+1}$ , the sum of values in  $F_{i+1}$ . However, if this sum  $\Sigma F_{i+1}$  is lower than  $th$ , expressed as a ratio (between 0 and 1) of the maximum among

all values  $\Sigma F_i$ , then the flux value at that point is set to NaN. By default when ‘Normal’ is toggled on, *th* is set to 0.03.

## PRESET MODEL

- *mirflux(..., ‘Buzzing’)* tailors the computation of spectral flux with the view to unveil salient aspects of sonic unsteadiness (Lartillot and Weisser, 2021). A spectrogram is first computed with emphasis on the auditory range, while filtering out frequency components of low energy. Only positive difference between successive frames is kept, in order to focus on increase of energy. Each flux value is finally normalized with respect to the summation of the spectrum over frequencies. This enables to highlight sudden increase in energy, while abstracting from the absolute amplitude of the overall energy. This corresponds to:

*mirflux(..., ‘Frame’, .05, .1, ‘Min’, 500, ‘Max’, 20000, ‘MinGate’, .01,  
‘NormalSpectrum’, ‘Global’, ‘Inc’, ‘Dist’, ‘City’)*

When the ‘Buzzing’ method is used for academic research, please cite the following publication:

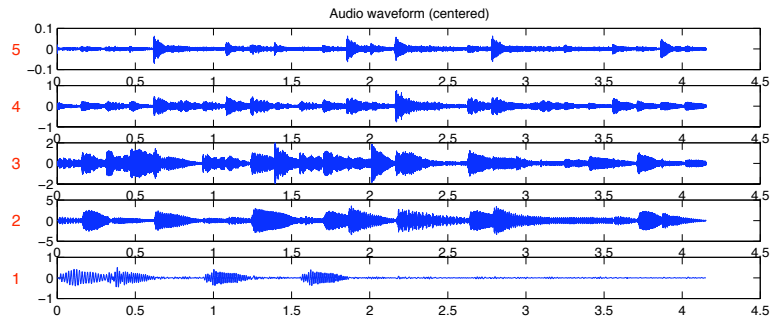
Lartillot, O. & Weisser, S. (2021). Roughness, Crackliness, Buzzingness, ...: Characterizations of Sonic Unsteadiness and Application to the Analysis of Traditional Music from Ethiopia, Kenya, Morocco and India. *16th International Conference on Music Perception and Cognition (ICMPC) & the 11th triennial conference of ESCOM*.

## *mirsum*

### SUMMATION OF FILTERBANK CHANNELS

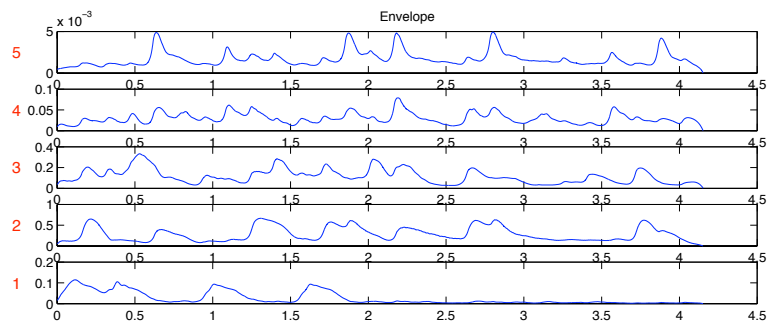
Once an audio waveform is decomposed into channels using a filterbank:

$$f = \text{mirfilterbank}(a)$$



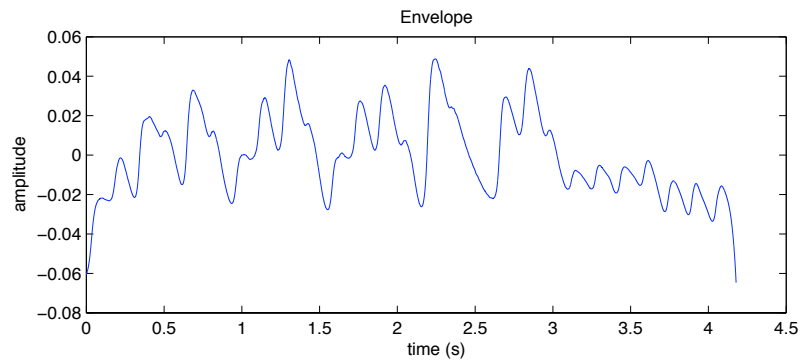
An envelope extraction, for instance, can be computed using this very simple syntax:

$$e = \text{mirenvelope}(f)$$



Then the channels can be summed back using the *mirsum* command:

$$s = \text{mirsum}(e)$$



The summation can be centered using the command:

$$s = \text{mirsum}(..., \textbf{Center})$$

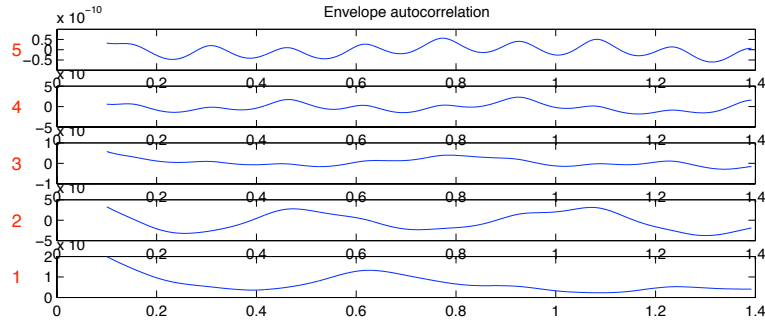
The summation can be divided by the number of channels using the command:

$$s = \text{mirsum}(..., \textbf{Mean})$$

## SUMMARY OF FILTERBANK CHANNELS

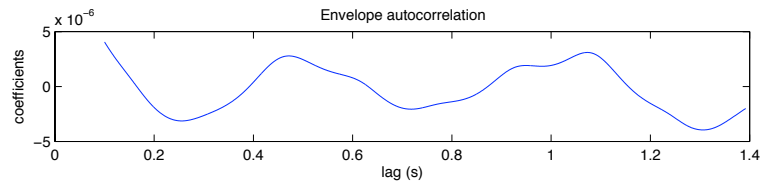
If we compute for instance an autocorrelation from the envelopes:

$$ac = \text{mirautocor}(e)$$



Then we can sum all the autocorrelation using exactly the same *mirsum* command:

$$s = \text{mirsum}(e)$$



This summation of non-temporal signals across channels is usually called *summary*.

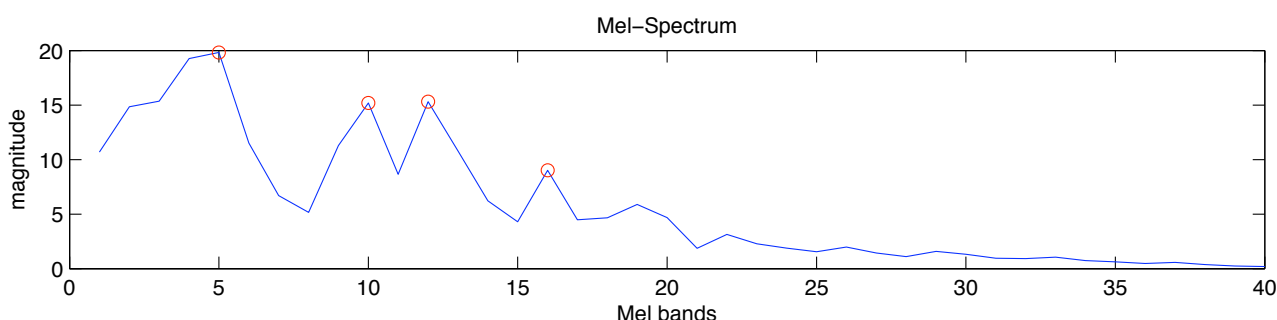
## *mirpeaks*

### PEAK PICKING

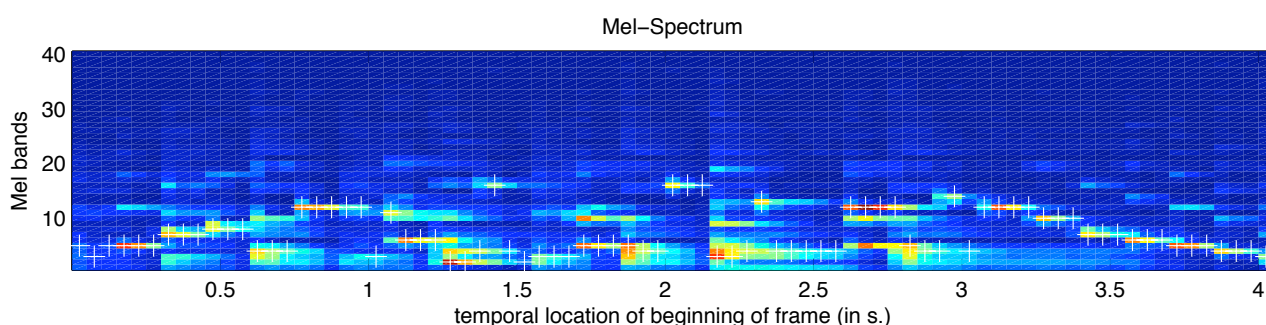
Peaks (or important local maxima) can be detected automatically from any data  $x$  produced in *MIRtoolbox* using the command

*mirpeaks*( $x$ )

If  $x$  is a curve, peaks are represented by red circles:



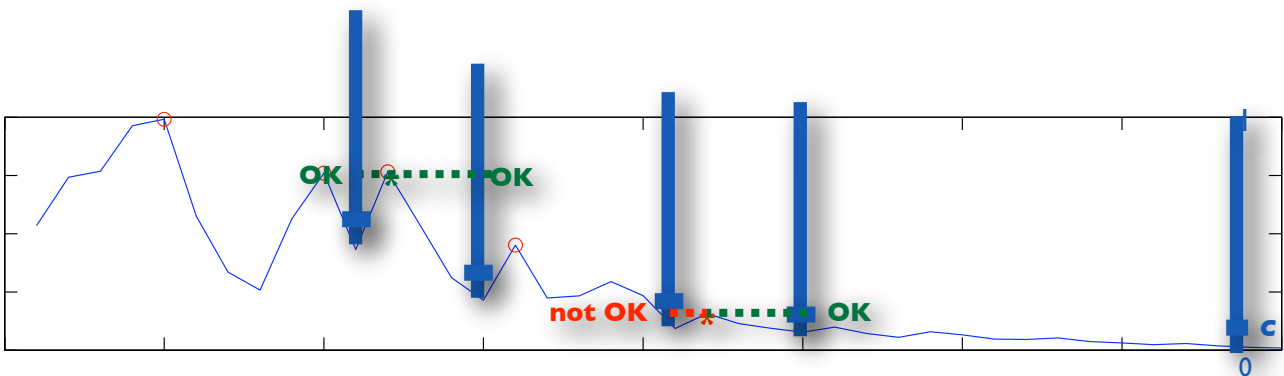
If  $x$  is a frame-decomposed matrix, peaks are represented by white crosses:



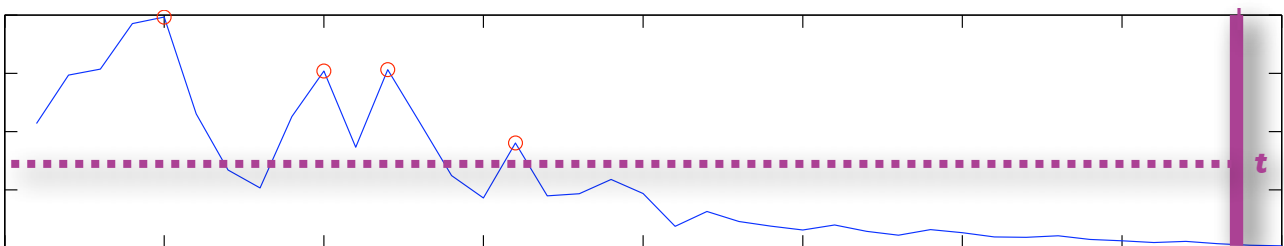
### PARAMETERS SPECIFICATION

- *mirpeaks*(..., **Total**,  $m$ ): only the  $m$  highest peaks are selected. If  $m = Inf$ , no limitation of number of peaks. Default value:  $m = Inf$
- Border effects can be specified:
  - *mirpeaks*(..., **NoBegin**) does not consider the first sample as a possible peak candidate. If a 'Contrast' parameter *cthr* is specified (cf. below), we will not consider either the first local maximum if the difference of amplitude with the beginning of the curve is below *cthr*.
  - *mirpeaks*(..., **NoEnd**) does not consider the last sample as a possible peak candidate. If a 'Contrast' parameter *cthr* is specified (cf. below), we will not consider either the last local maximum if the difference of amplitude with the end of the curve is below *cthr*.

- *mirpeaks*(..., '**Order**', *o*) specifies the ordering of the peaks.
  - *o* = '**Amplitude**' orders the peaks from highest to lowest (Default choice.)
  - *o* = '**Abscissa**' orders the peaks along the abscissa axis.
- *mirpeaks*(..., '**Valleys**') detect valleys (local minima) instead of peaks.
- *mirpeaks*(..., '**Contrast**', *cthr*): A given local maximum will be considered as a peak if the difference of amplitude with respect to both the previous and successive local minima (when they exist) is higher than the threshold *cthr*. This distance is expressed with respect to the total amplitude of the input signal: a distance of 1, for instance, is equivalent to the distance between the maximum and the minimum of the input signal. Default value: *cthr* = 0.1



- *mirpeaks*(..., '**SelectFirst**', *fthr*): If the '**Contrast**' selection has been chosen, this additional option specifies that when one peak has to be chosen out of two candidates, and if the difference of their amplitude is below the threshold *fthr*, then the most ancient one is selected. Option toggled off by default. Default value if toggled on: *fthr* = *cthr*/2
- *mirpeaks*(..., '**Threshold**', *thr*): A given local maximum will be considered as a peak if its normalized amplitude is higher than this threshold *thr*. A given local minimum will be considered as a valley if its normalized amplitude is lower than this threshold. The normalized amplitude can have value between 0 (the minimum of the signal in each frame) and 1 (the maximum in each frame). Default value: *thr*=0 for peaks, *thr* = 1 for valleys.



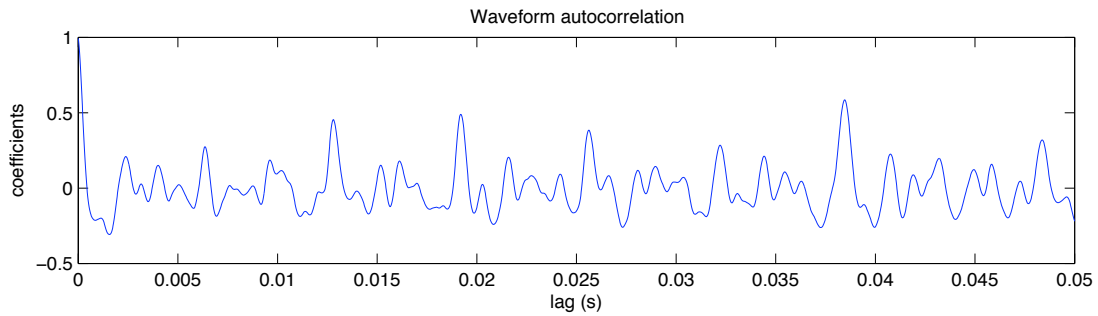


- *mirpeaks(..., 'Interpol', i)* estimates more precisely the peak position and amplitude using interpolation. Performed only on data with numerical abscissae axis.
  - *i* = "", 'no', 'off', 0: no interpolation
  - *i* = 'Quadratic': quadratic interpolation. (default value).
- *mirpeaks(..., 'Reso', r)* removes peaks whose abscissa distance to one or several higher peaks is lower than a given threshold. *r* is the numerical value of the threshold, or it can be 'Semi-Tone', corresponding to a ratio between the two peak frequency positions equal to  $2^{(1/12)}$ . By default, out of two conflicting peaks, the higher peak remains. If the keyword 'First' is added, the peak with lower abscissa value remains instead.
- *mirpeaks(..., 'Pref', c, std)* indicates a region of preference for the peak picking, centered on the abscissa value *c*, with a standard deviation of *std*.
- *mirpeaks(..., 'Nearest', t, s)* takes the peak nearest a given abscissa values *t*. The distance is computed either on a linear scale (*s* = 'Lin') or logarithmic scale (*s* = 'Log'). When using the 'Nearest' option, only one peak is extracted.

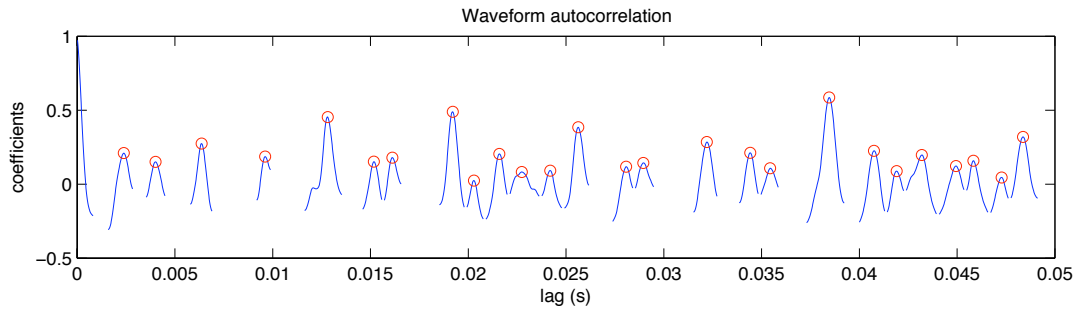
The 'Total' parameter can then be used to indicate the number of peaks to preselect before the 'Nearest' selection. If 'Total' was still set to 1, it is then ignored – i.e., forced to *Inf* – in order to preselect all possible peaks.

- *mirpeaks(..., 'Normalize', n)* specifies whether frames are normalized globally or individually.
  - *n* = 'Global' normalizes across the whole frames (as well as across the whole segments) of each audio file from 0 to 1 (default choice).
  - *n* = 'Local': normalizes each segment, each frame, from 0 to 1 separately.
  - *n* = 'No' or 0: no normalization at all.
- *mirpeaks(..., 'Normalize', 'Local', 'LocalFactor', f)* normalizes each frame individually but also based on a certain number of previous frames. More precisely, the normalisation  $N(i)$  for a given frame *i* is the maximum between the maximum in the given frame *i* and *f* times the normalisation  $N(i-1)$  for the previous frame. When 'LocalFactor' is used, *f* is set by default to .99. This enables to have an adaptive normalisation based on the local context, but in the same time avoid focusing too much on abrupt silences.
- *mirpeaks(..., 'Extract')* extracts from the curves all the positive continuous segments (or "curve portions") where peaks are located. First, a low-pass filtered version of the curve is computed, on which the temporal span of the positive lobes containing each peak are stored. The output consists of the part of the original non-filtered curve corresponding to the same temporal span. For instance:

*ac = mirautocor('ragtime.wav')*

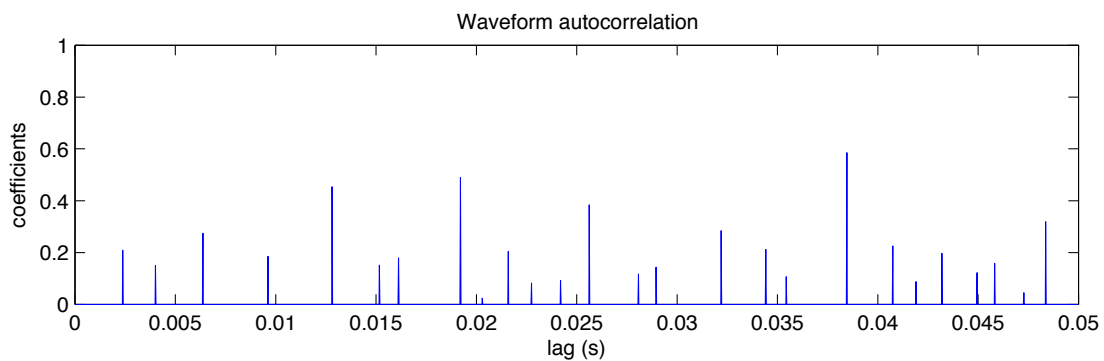


*mirpeaks(ac, 'Extract')*



- *mirpeaks(..., 'Only')*, keeps from the original curve only the data corresponding to the peaks, and zeroes the remaining data.

*mirpeaks(ac, 'Only')*

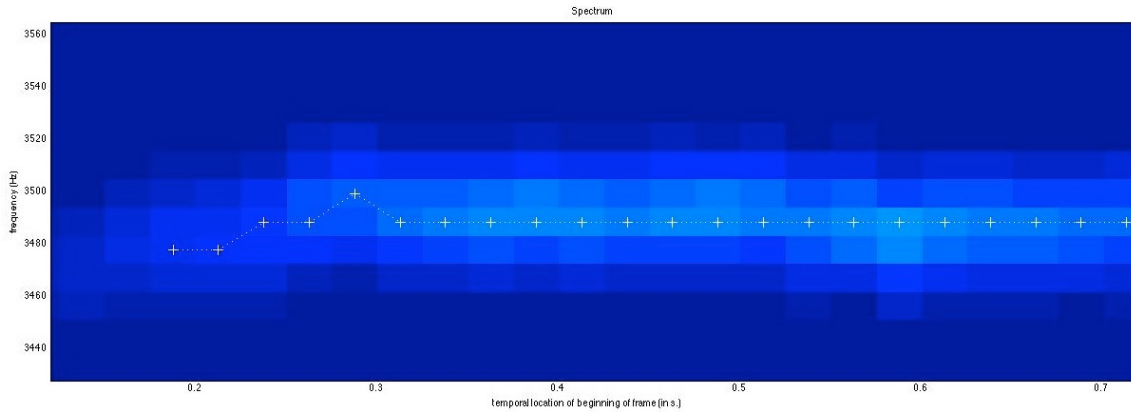


- *mirpeaks(..., 'Track', t)*, where the input is some frame-decomposed vectorial data – such as spectrogram, for instance –, tracks peaks along time using McAulay & Quatieri's (1986) method: lines are drawn between contiguous temporal series of peaks that are sufficiently aligned. If a value *t* is specified, the variation between successive frames is tolerated up to *t*, *ex-*

*pressed using the abscissae unit*. For instance, the figure below shows the result (zoomed) of the following commands:

```
s = mirspectrum('trumpet', 'Frame');
```

```
mirpeaks(s, 'Track', 25)
```



- *mirpeaks(..., 'CollapseTracks', t)*, collapses tracks into one single track, and remove small track transitions, of length shorter than *ct* samples. Default value: *ct* = 7.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **'PeakPos'**: the abscissae position of the detected peaks, in sample index,
- **'PeakPosUnit'**: the abscissae position of the detected peaks, in the default abscissae representation. If the input is a *mirscalar* object (i.e., where there is one numerical value associated to each successive frame, such as *mirflux*, *mirnovelty*, etc.), this returns the frame positions of the detected peaks, each frame position being the start and end time of the frame in a column vector.
- **'PeakPrecisePos'**: a more precise estimation of the abscissae position of the detected peaks computed through interpolation, in the default abscissae representation,
- **'PeakVal'**: the ordinate values associated to the detected peaks,
- **'PeakPreciseVal'**: a more precise estimation of the ordinate values associated to the detected peaks, computed through interpolation,
- **'PeakMode'**: the mode values associated to the detected peaks,

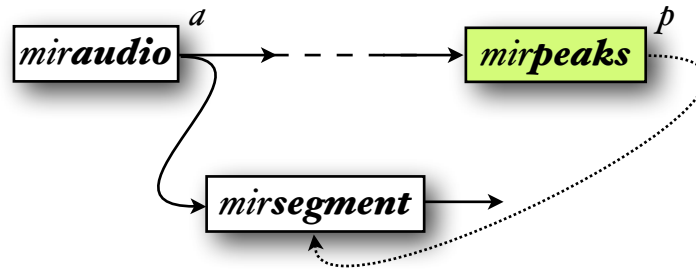
- ***TrackPos***: the abscissae position of the peak tracks, in sample index,
- ***TrackPosUnit***: the abscissae position of the peak tracks, in the default abscissae representation,
- ***TrackVal***: the ordinate values of the peak tracks..

## *mirsegment*

### SEGMENTATION

- An audio waveform  $a$  can be segmented using the output  $p$  of a peak picking from data resulting from  $a$  itself, using the following syntax:

$$sg = \text{mirsegment}(a, p)$$



If  $p$  is a frame-decomposed scalar curve, the audio waveform  $a$  will be segmented at the middle of each frame containing a peak.

- An audio waveform  $a$  can be segmented at the beginning of the pitch events detected using *mirpitch*:

$$p = \text{mirpitch}(a, \text{'Frame'}, \text{'Segment'}, \dots)$$

$$sg = \text{mirsegment}(a, p)$$

- An audio waveform  $a$  can also be segmented manually, based on temporal position directly given by the user, in the form:

$$sg = \text{mirsegment}(a, v)$$

where  $v$  is an array of numbers corresponding to time positions in seconds.

If  $a$  is a set of  $N$  audio waveforms, and if  $v$  is a matrix of  $N$  columns, then each column  $i$  indicates the segmentation points for the waveform  $i$ . If  $v$  has only one column, then all waveforms use that same column of segmentation points.

- A segmented audio waveform  $a$  can be further segmented manually, based on temporal position directly given by the user, in the form:

$$sg = \text{mirsegment}(a, v)$$

where  $v$  is a cell array, where each cell, corresponding to one audio file, is itself a cell array, where each cell, corresponding to one segment, is an array of numbers corresponding to time positions in seconds.

- Automated segmentation methods are provided as well, that can be called using the syntax:

$$sg = \text{mirsegment}(a, m)$$

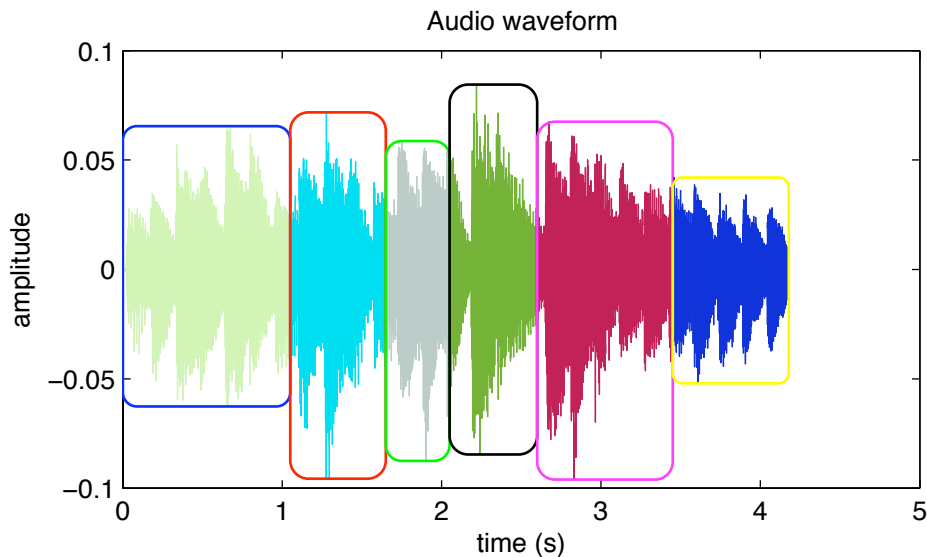
where  $m$  is the name of one of the following segmentation methods: ‘Novelty’ (default, cf. *mirnovelty*), ‘HCDF’ (cf. *mirhcdf*) or ‘RMS’ (cf. *mirrms*).

*mirsegment* accepts uniquely as main input a ***miraudio*** objects not frame-decomposed, not channel decomposed, and not already segmented. Alternatively, **file name** or the ‘**Folder**’ keyword can be used as well.

The first argument of the *mirsegment* function is the audio file that needs to be segmented. It is possible for instance to compute the segmentation curve using a downsampled version of the signal and to perform the actual segmentation using the original audio file.

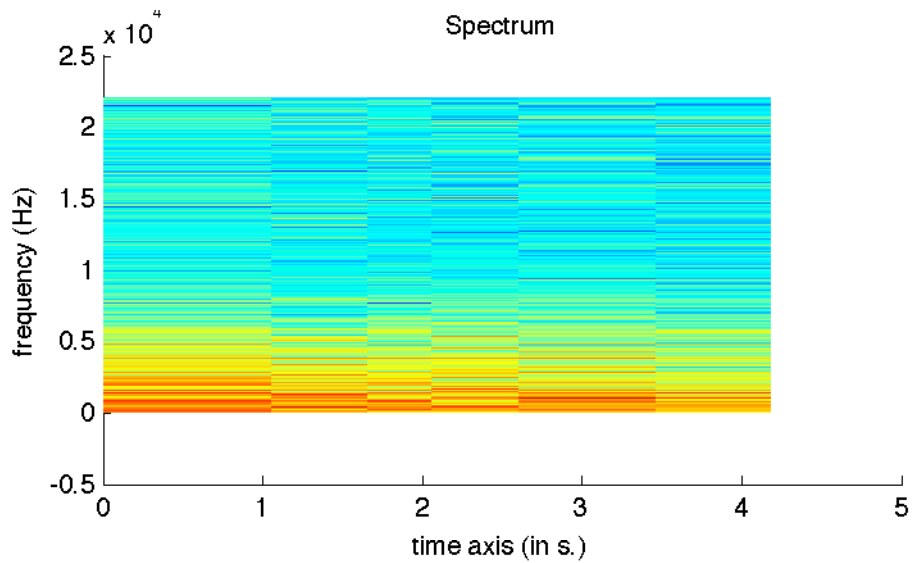
## EXAMPLE

$$sg = \text{mirsegment}('ragtime.wav', 'Novelty', 'KernelSize', 32)$$



The output can be sent to any further analysis, for instance:

$$sp = \text{mirspectrum}(sg, 'dB')$$



## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method.

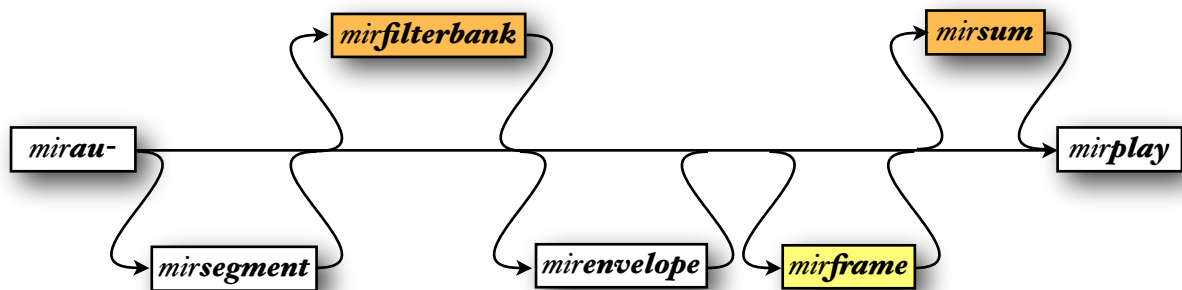
- ***FramePos***: For segmented data, this returns a cell array where each cell is a matrix containing the starting and ending temporal positions of the frames in each successive segment. When there is no frame decomposition, each cell contains simply the starting and ending time of each successive segment.

## *mirplay*

### SONIFICATION OF THE RESULT

Certain classes of temporal data can be sonified:

- ***miraudio*** objects: the waveform is directly played, and
  - if the audio waveform is segmented (using ***mirsegment***), segments are played successively with a short burst of noise in-between;
  - if the audio waveform is decomposed into channels (using ***mirfilterbank***), channels are played successively from low to high register;
  - if the audio is decomposed into frames (using ***mirframe***), frames are played successively;
- **file name** or the **'Folder'** keyword: same behavior than for *miraudio* objects;
- ***mirenvelope*** objects (frame-decomposed or not) are sonified using a white noise modulated in amplitude by the envelope, and
  - if peaks have been picked on the envelope curve (using ***mirpeaks***), they are sonified using a short impulse sound;
- ***mirpitch*** results: each extracted frequency is sonified using a sinusoid.



### FRAME DECOMPOSITION

*mirplay*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

### OPTIONS

- *mirplay*(..., **'Channel'**, *i*) plays the channel(s) of rank(s) indicated by the array *i*.
- *mirplay*(..., **'Segment'**, *k*) plays the segment(s) of rank(s) indicated by the array *k*.



- *mirplay*(..., '**Sequence**', *l*) plays a sequence of audio files using the order indicated by the array *l*.
- *mirplay*(..., '**Increasing**', *d*) plays the sequences in increasing order of *d*, which can be either an array of number or a *mirscalar* data (i.e., a scalar data returned by *MIRtoolbox*).
- *mirplay*(..., '**Decreasing**', *d*) plays the sequences in decreasing order of *d*, which can be either an array of number or a *mirscalar* data (i.e., a scalar data returned by *MIRtoolbox*).
- *mirplay*(..., '**Every**', *s*) plays every *s* sequence, where *s* is a number indicating the step between sequences.
- *mirplay*(..., '**Burst**', 'No') toggles off the burst sound between segments.

Example:

```
e = mirenvelope('Folder');
```

```
rms = mirrms('Folder');
```

```
mirplay(e, 'increasing', rms, 'every', 5)
```

## *mirplayer*

### GRAPHICAL PLAYER INTERFACE

*mirplayer(features)* opens a GUI showing the graphical representation of the *MIRtoolbox* object *features*. *features* is a mirtoolbox object such as *mirstruct* or *mirscalar*, or a nested MATLAB structure of mirtoolbox objects.

*mirplayer(features, select)* displays only songs specified by an array of song indices *select*.

### EXAMPLES

The examples assume that the folder *MIRToolboxDemos* is in the MATLAB path.

```
b = mirbrightness('ragtime.wav','Frame');
```

```
mirplayer(b)
```

```
features.brightness = mirbrightness('ragtime.wav','Frame');
```

```
features.filterbank = mirfilterbank('ragtime.wav');
```

```
mirplayer(features)
```

```
features = mirfeatures('ragtime.wav');
```

```
mirplayer(features)
```

```
a = miraudio('ragtime.wav');
```

```
mirplayer(a)
```

Display only the first and the fifth audio file:

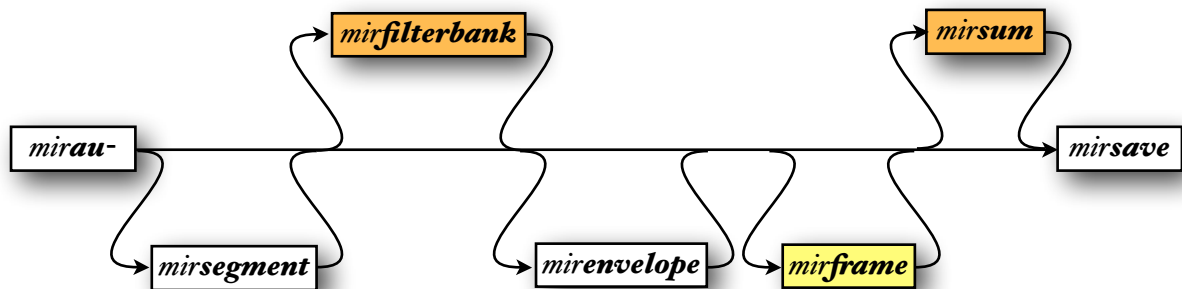
```
mirplayer(features,[1,5])
```

## *mirsave*

### SAVING AUDIO RENDERING INTO FILES

Certain classes of temporal data can be saved:

- ***miraudio*** objects: the waveform is directly saved, and
  - if the audio waveform is segmented (using ***mirsegment***), segments are concatenated with a short burst of noise in-between;
  - if the audio waveform is decomposed into channels (using ***mirfilterbank***), each channel is saved in a separate file;
  - if the audio is decomposed into frames (using ***mirframe***), frames are concatenated;
- **file name** or the **'Folder'** keyword: same behavior than for *miraudio* objects;
- ***mirenvelope*** objects (frame-decomposed or not) are sonified using a white noise modulated in amplitude by the envelope,
- ***mirpitch*** results: each extracted frequency is sonified using a sinusoid.



### FRAME DECOMPOSITION

*mirsave*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

### OPTIONS

- The name and extension of the saved file can be specified in different ways, as shown in the tables below.
  - By default, the files are saved in WAV format, using the extension *'mir.sav'* in order to lower the risk of overwriting original audio files.

- If the string **'au'** is indicated as second argument of *mirsave*, the audio will be saved in AU format.
- A string can be indicated as second argument of *mirsave*.
  - If the *miraudio* object to be saved contains only one audio file, the specified string will be used as the name of the new audio file.
  - If the *miraudio* object to be saved contains several audio files, the specified string will be concatenated to the original name of each audio file.
- If the second argument of *mirsave* is a string ended by **'au'**, the file name will follow the convention explained in the previous point, and the files will be saved in AU format.

<i>a = miraudio('mysong.au')</i>		<i>mysong.au</i>
<i>mirsave(a)</i>		<i>mysong.mir.wav</i>
<i>mirsave(a,'new')</i>		<i>new.wav</i>
<i>mirsave(a,'.au')</i>		<i>mysong.mir.au</i>
<i>mirsave(a,'new.au')</i>		<i>new.au</i>
<i>a = miraudio('Folder')</i>	<i>song1.wav</i>	<i>song2</i>
<i>mirsave(a)</i>	<i>song1.mir.wav</i>	<i>song2.m</i>
<i>mirsave(a,'new')</i>	<i>song1new.wav</i>	<i>song2ne</i>
<i>mirsave(a,'.au')</i>	<i>song1.mir.au</i>	<i>song2.n</i>
<i>mirsave(a,'new.au')</i>	<i>song1new.au</i>	<i>song2n</i>

*Diverse ways of saving into an audio file*

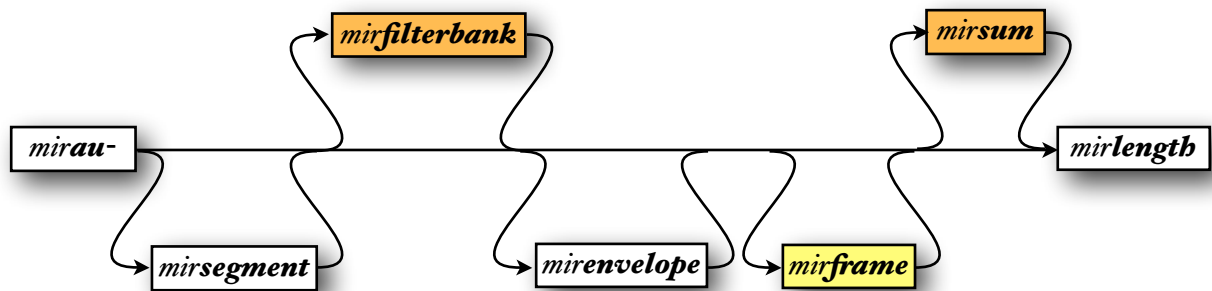
*Diverse ways of saving as a batch of audio files*

- *mirsave(a, filename, 'SeparateChannels')* save each separate channel in a different file. The channel number is added to the file name, before any file extension.
- *mirsave(a, filename, 'SeparateSegments')* save each separate segment in a different file. The segment number is added to the file name, before any file extension.

## *mirlength*

### TEMPORAL LENGTH OF SEQUENCES

*mirlength* returns the temporal length of the temporal sequence given in input, which can be either an audio waveform (*miraudio*) or an envelope curve (*mirenvelope*). If the input was decomposed into segments (*mirsegment*), *mirlength* returns a curve indicating the series of temporal duration associated with each successive segment.



### OPTIONS

- *mirlength*(..., '**Unit**', *u*) specifies the length unit. The possible values are:
  - *u* = 'Second': duration in seconds (Default choice).
  - *u* = 'Sample': length in number of samples.

## FILTERBANK DECOMPOSITION WITH FILTER-DEPENDENT WINDOWING

*Mirfdw* computes spectral sub-band features using Filter Dependent Windowing (FDW).

Similar to wavelet transform, FDW uses a separate frame size based on the central frequency of each sub-band filter. So the following function is used:

$$FDTF = \frac{100}{Sub - band\ central\ freq}$$

$$Sub - band\ central\ freq = \begin{cases} \sqrt{(high\ cutoff\ frequency * low\ cutoff\ frequency)} , & \text{if } \frac{high\ cutoff\ frequency}{low\ cutoff\ frequency} \geq 1.1 \\ \frac{1}{2}(high\ cutoff\ frequency + low\ cutoff\ frequency), & \text{if } \frac{igh\ cutoff\ frequency}{low\ cutoff\ frequency} < 1.1 \end{cases}$$

So when we apply this simple function we get the following sub-band windows:

Sub-band Filter	Frequency Range (Hertz)	Central Frequency (Hertz)	Frame size (Seconds)
No 1	0 - 50	25	4
No 2	50 - 100	71	2,92
No 3	100 - 200	141	4,91
No 4	200 - 400	283	2,46
No 5	400 - 800	566	1,23
No 6	800 - 1600	1131	0,61
No 7	1600 - 3200	2263	0,31
No 8	3200 - 6400	4525	0,15
No 9	6400 - 12800	9051	0,08
No 10	12800 - 22050	16800	0,06

Initial algorithm developed by Fabi Prezja for his Master's thesis<sup>3</sup>, University of Jyväskylä, 2018.

## SYNTAX

*mirfdw(filename, method, order)* where

- *filename* is the name of the audio file used as input

<sup>3</sup> <https://jyx.jyu.fi/handle/123456789/60963#>

- *method* is the method applied to each channel after the filterbank decomposition. By default, *method* = *@mirflux*. (Somewhat related to *mirflux*(..., 'SubBand')). Other suggested methods: *@mirskewness*, *@mirzerocross*, *@mirkurtosis*, *@mirentropy*.
- *order* is the filter order used as option in *mirfilterbank*(..., 'Manual'). Default value: *order* = 4

### 3. FEATURE EXTRACTORS

The musical feature extractors can be organized along main musical dimensions: dynamics, rhythm, timbre, pitch and tonality.

#### 3.1. Dynamics

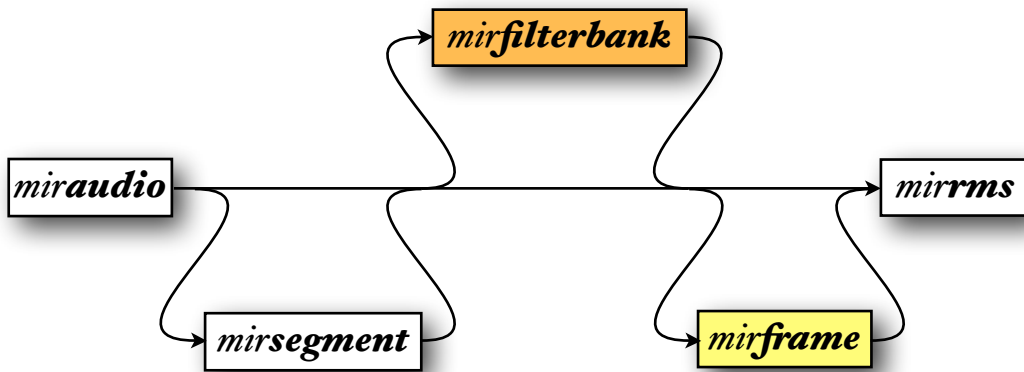
##### *mirrms*

##### ROOT-MEAN-SQUARE ENERGY

The global energy of the signal  $x$  can be computed simply by taking the root average of the square of the amplitude, also called root-mean-square (RMS):

$$x_{\text{rms}} = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2} = \sqrt{\frac{x_1^2 + x_2^2 + \dots + x_n^2}{n}}$$

##### FLOWCHART INTERCONNECTIONS



*mirrms* accepts as input data type either:

- ***miraudio*** objects, where the audio waveform can be segmented (using ***mirsegment***), decomposed into channels (using ***mirfilterbank***), and/or decomposed into frames (using ***mirframe***),
- **file name** or the **'Folder'** keyword,
- other vectorial objects, such as *mirspectrum*, are accepted as well, although a warning is displayed in case this particular configuration was unintended, i.e., due to an erroneous script.

The following command orders the computation of the RMS related to a given audio file:



*mirrms('ragtime.wav')*

which produce the resulting message in the *Command Window*:

The RMS energy related to file ragtime.wav is 0.017932

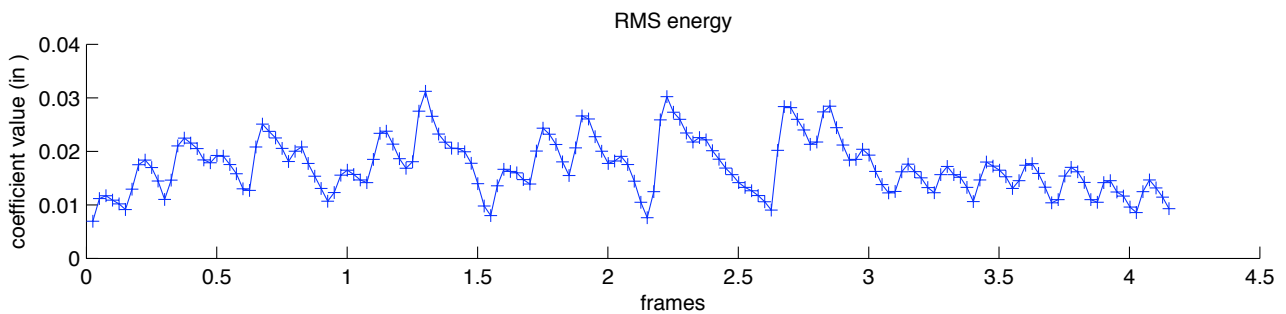
## FRAME DECOMPOSITION

*mirrms(..., 'Frame', ...)* performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. 'Frame' section.

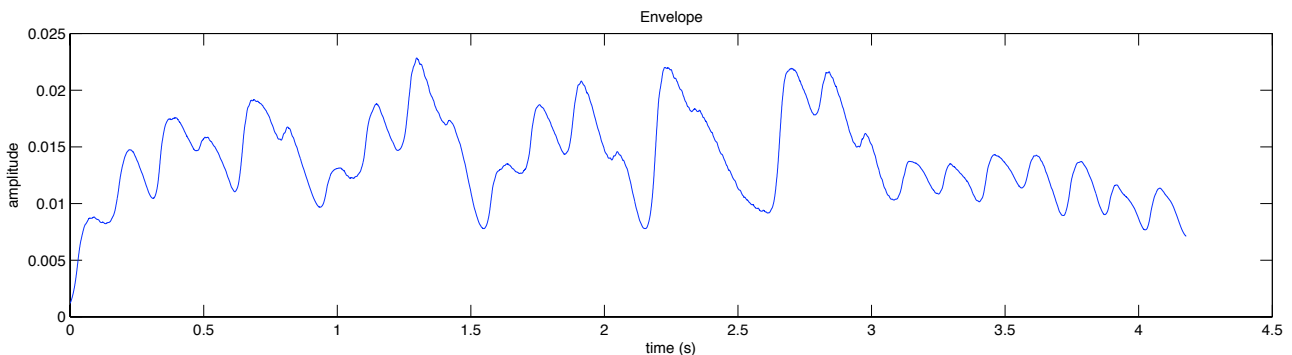
For instance:

*mirrms('ragtime.wav', 'Frame')*

we obtain a temporal evolution of the energy:



We can note that this energy curve is very close to the envelope:



## OPTIONS

- *mirrms(..., 'Median')* replaces the mean by the median, corresponding to what could be called a “root-median-square”. If 'Frame' is not used, chunk decomposition does not work and should therefore be avoided using *mirchunklim(Inf)*.

- *mirrms*(..., '**Warning**', 0) toggles off the warning message, explained above, i.e., when the input is a vectorial data that is not a *miraudio* object.

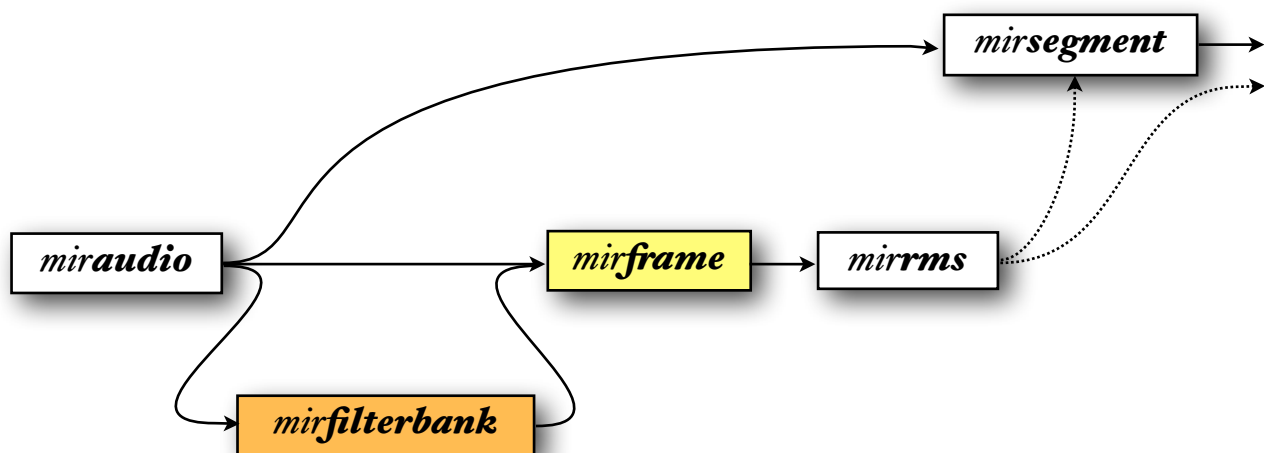
## *mirsegment(..., 'RMS')*

Segmentation at positions of long silences. A frame decomposed RMS is computed using *mirrms* (with default options), and segments are selected from temporal positions where the RMS rises to a given 'On' threshold, until temporal positions where the RMS drops back to a given 'Off' threshold.

### OPTIONS

- *mirsegment(..., 'Off',  $t_1$ )* specifies the RMS 'Off' threshold. Default value:  $t_1 = .01$
- *mirsegment(..., 'On',  $t_2$ )* specifies the RMS 'On' threshold. Default value:  $t_2 = .02$

### FLOWCHART INTERCONNECTIONS



*mirsegment* accepts uniquely as main input a *miraudio* objects not frame-decomposed, not channel decomposed, and not already segmented. Alternatively, **file name** or the '**Folder**' key-word can be used as well.

*mirsegment(..., 'RMS')* can return several outputs:

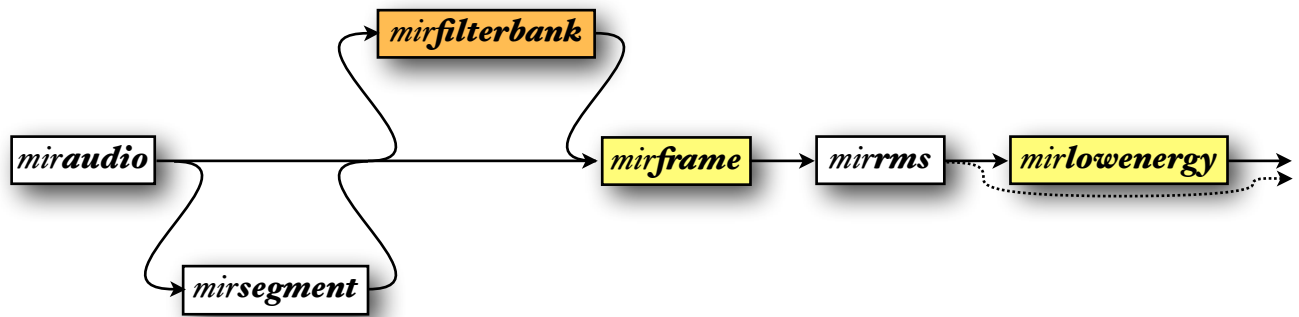
1. the segmented audio waveform itself,
2. the RMS curve (*mirrms*).

## *mir*lowenergy

### DESCRIPTION

The energy curve can be used to get an assessment of the temporal distribution of energy, in order to see if it remains constant throughout the signal, or if some frames are more contrastive than others. One way to estimate this consists in computing the low energy rate, i.e. the percentage of frames showing less-than-average energy (Tzanetakis and Cook, 2002).

### FLOWCHART INTERCONNECTIONS



*mirlowenergy* accepts as input data type either:

- ***mirrms* frame-decomposed** data,
- ***miraudio*** objects, where the audio waveform can be segmented (using *mirsegment*), decomposed into channels (using *mirfilterbank*). The audio waveform is decomposed into frames if it was not decomposed yet.
- **file name** or the **'Folder'** keyword: same behavior than for *miraudio* objects.

*mirlowenergy* can return several outputs:

1. the low-energy rate itself and
2. the ***mirrms* frame-decomposed** data.

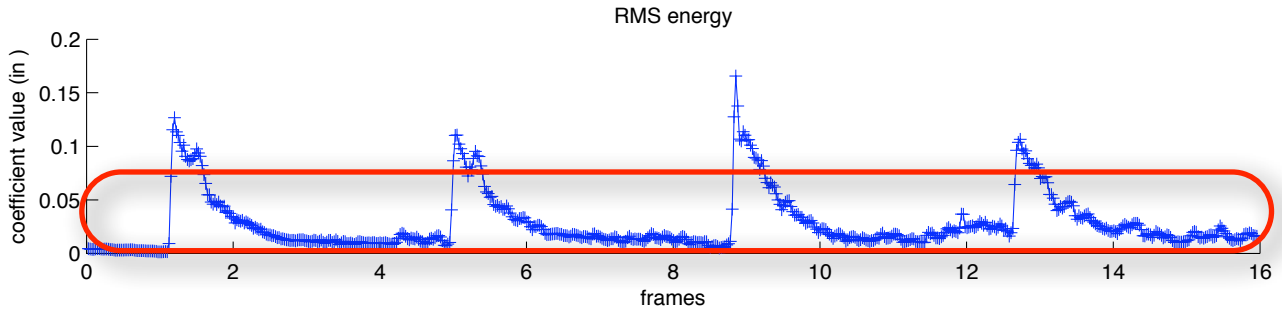
### FRAME DECOMPOSITION

*mirlowenergy*(..., **'Frame'**, ...) specifies the frame configuration, with by default a frame length of 50 ms and half overlapping. For the syntax, cf. the previous *mirframe* vs. 'Frame' section.

### EXAMPLES

If we take for instance this energy curve:

$$rI = \text{mirrms}(aI, \text{'Frame'})$$



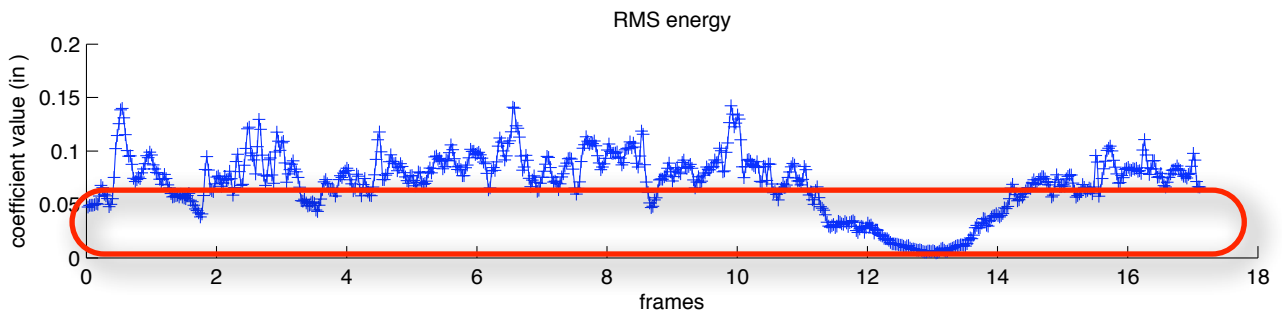
We can see that due to some rare frames containing particularly high energy, most of the frames are below the average RMS. And indeed if we compute the low-energy rate

$$\text{mirlowenergy}(r1)$$

we obtain the value 0.71317.

For this opposite example:

$$r2 = \text{mirrms}(a2, 'Frame')$$



there are two kind of frames basically, those that have quite constant high energy, and fewer that have very low energy. Hence most of the frames are over the average energy, leading to a low low-energy rate:

$$\text{mirlowenergy}(r2)$$

equal to 0.42398.

## OPTIONS

- *mirlowenergy*(..., '**Threshold**', *t*) expressed as a ratio to the average energy over the frames.  
Default value:  $t = 1$
- *mirlowenergy*(..., '**ASR**') computes the Average Silence Ratio (Feng, Zhuang, Pan, 2003), which corresponds to a RMS without the square-root, and a default threshold set to  $t = .5$

### 3.2. Rhythm

The estimation of rhythmicity in the audio signal can be performed using the basic operators we introduced previously.

#### *mirfluctuation*

#### RHYTHMIC PERIODICITY ALONG AUDITORY CHANNELS

One way of estimating the rhythmic is based on spectrogram computation transformed by auditory modeling and then a spectrum estimation in each band (Pampalk et al., 2002).

The implementation proposed in *MIRtoolbox* includes a subset of the series of operations proposed in Pampalk et al.:

1. First a power spectrogram is computed

- on frames of 23 ms with a hop rate *hr*, set by default to 80 Hz, but is automatically raised to the double of the maximum fluctuation frequency range *m*, used in step 2.
- The Terhardt outer ear modeling is computed.
- A multi-band redistribution of the energy is performed along the '**Bark**' bands decomposition.
  - *mirfluctuation*(..., '**Mel**') performs a decomposition into 40 Mel bands instead of 20 Bark bands.
- Masking effects are estimated on the multi-band distribution.
- Finally the amplitudes are represented in dB scale.

This is summarized in one *MIRtoolbox* command line:

$$s = \textit{mirspectrum}(\dots, \textit{Frame}, .023, \textit{s}, \textit{hr}, \textit{Hz}, \dots$$
$$\textit{Power}, \textit{Terhardt}, \textit{b}, \textit{Mask}, \textit{dB})$$

where *b* is either '*Bark*' or '*Mel*'.

2. Then a FFT is computed on each band:

- Frequencies range from 0 to  $m$  Hz, where  $m$  is set by default to 10 Hz can be controlled by a '**Max**' option: `mirfluctuation(..., 'Max', m)`
- The frequency resolution of the FFT  $mr$  is set by default to .01 Hz and can also be controlled by a '**MinRes**' option: `mirfluctuation(..., 'MinRes', mr)`
- The amplitude modulation coefficients are weighted based on the psychoacoustic model of the fluctuation strength (Fastl, 1982).

This is summarized in one *MIRtoolbox* command line:

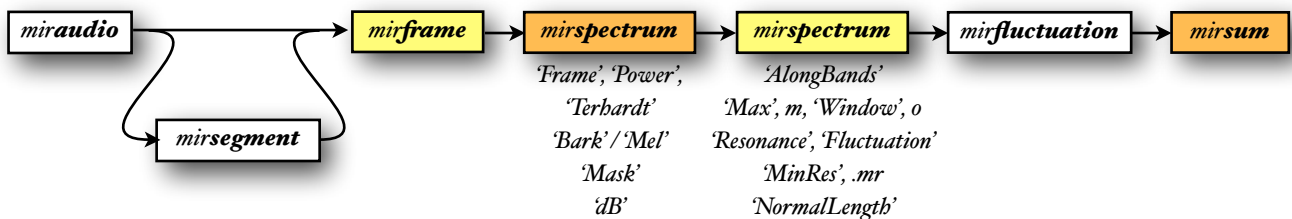
$$f = \text{mirspectrum}(s, 'AlongBands', 'Max', m, 'MinRes', mr, 'Window', o, \dots \\ 'Resonance', 'Fluctuation', 'NormalLength')$$

We can see in the matrix the rhythmic periodicities for each different Bark band.

3. `mirfluctuation(..., 'Summary')` subsequently sums the resulting spectrum across bands, leading to a spectrum summary, showing the global repartition of rhythmic periodicities:

$$\text{mirsum}(f)$$

## FLOWCHART INTERCONNECTIONS



`mirfluctuation` accepts as input data type either:

- **mirspectrum frame-decomposed** objects (i.e., spectrograms),
- **miraudio** objects, where the audio waveform can be segmented (using `mirsegment`). The audio waveform is decomposed into frames if it was not decomposed yet.
- **file name** or the '**Folder**' keyword: same behavior than for `miraudio` objects.

## FRAME DECOMPOSITION

- *mirfluctuation*(..., '**InnerFrame**',  $l$ ,  $r$ ) specifies the spectrogram frame length  $l$  (in second), and, optionally, the frame rate  $r$  (in Hertz), with by default a frame length of 23 ms and a frame rate of 80 Hz.
- *mirfluctuation*(..., '**Frame**',  $l$ ,  $r$ ) computes fluctuation using a window moving along the spectrogram, whose length  $l$  (in second) and frame rate  $r$  (in Hertz) can be specified as well, with by default a frame length of 1 s and a frame rate of 10 Hz.

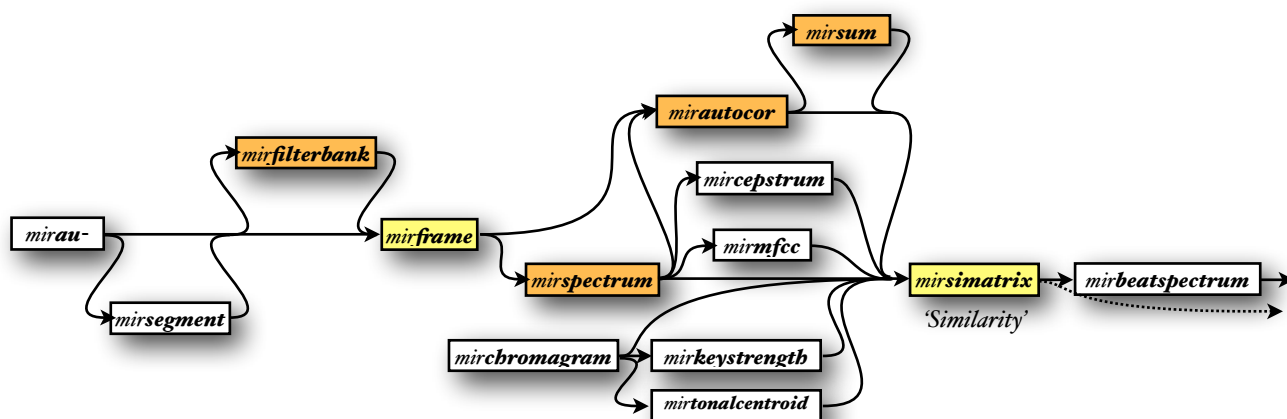


## *mirbeatspectrum*

### BEAT SPECTRUM

The beat spectrum has been proposed as a measure of acoustic self-similarity as a function of time lag, and is computed from the similarity matrix (cf. *mirsimatrix*) (Foote, Cooper and Nam, 2002).

### FLOWCHART INTERCONNECTIONS



One parameter related to *mirsimatrix* is accessible in *mirbeatspectrum*:

- **'Distance'**.

*mirbeatspectrum* accepts either:

- *mirsimatrix* objects,
- *mirspectrum* frame-decomposed objects,
- *miraudio* objects: in this case, the similarity matrix will be based on the mfcc (*mirmfcc*), computed from ranks 8 to 33. The audio waveform is decomposed into frames if it was not decomposed yet
- **file name** or the **'Folder'** keyword: same behavior as for *miraudio* objects,
- other frame-decomposed analysis.

*mirbeatspectrum* can return several outputs:

1. the beat spectrum curve itself, and
2. the similarity matrix (*mirsimatrix*).

## FRAME DECOMPOSITION

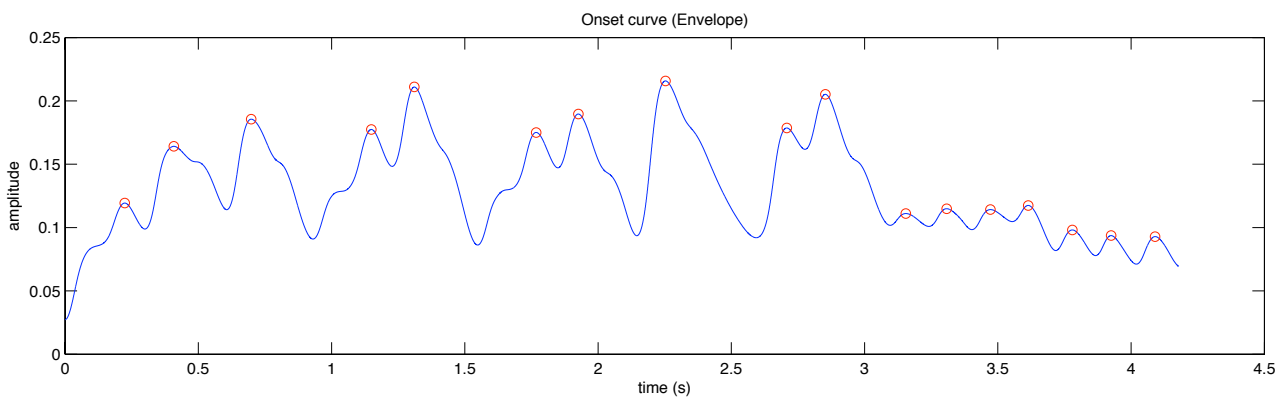
*mirbeatspectrum*(..., ***Frame***, ...) specifies the frame configuration, with by default a frame length of 25 ms and a hop factor of 10 ms. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *Frame* section.

## *miреvents* (previously *miронsets*)

### TEMPORAL LOCATION OF EVENTS

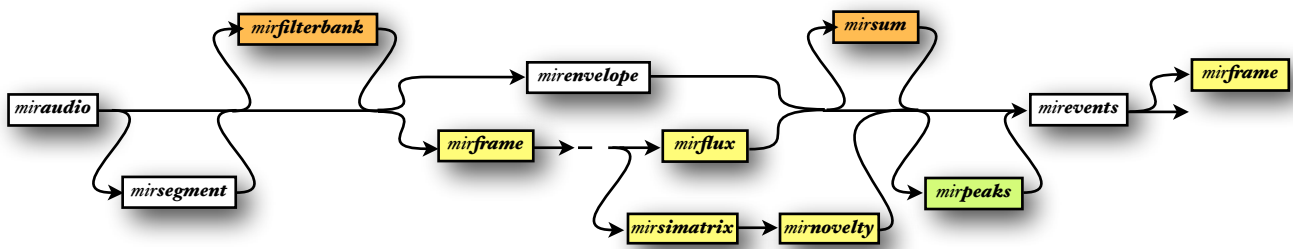
Detection curve, showing the successive bursts of energy corresponding to the successive pulses or sound events. A peak picking is automatically performed on the detection curve, in order to show the estimated positions of the events, such as notes, chords, etc.

Related to event detection is the detection of attack and decay phases (cf. below). It is also possible to detect attack and decay phases without computing a detection curve, using the ‘Waveform’ method.



*miреvents*('ragtime.wav')

### FLOWCHART INTERCONNECTIONS



The detection curve can be computed in various ways:

- *miреvents*(..., '**Envelope**') computes an amplitude envelope, using *mirenvelope* (default choice). The envelope extraction can be specified, as in *mirenvelope*:
  - either the '**Spectro**' option (default):
    - *miреvents*(..., '**SpectroFrame**', *fl*, *fb*) species the frame length *fl* (in s.) and the hop factor *fb* (as a value between 0 and 1). Default values: *fl* = .1 s., *fb* = .1

- the frequency reassignment method can be specified: '**Freq**' (default), '**Mel**', '**Bark**' or '**Cents**' (cf. *mirspectrum*).
- *mirevents*(..., '**PowerSpectrum**', 0) turns off the computation of the power of the spectrum. When the '**Attacks**' and/or '**Decays**' options are used (cf. below), '**PowerSpectrum**' is toggled off by default.
- *mirevents*(..., '**Terhardt**') toggles on the '**Terhardt**' operation (cf. *mirspectrum*).
- *mirevents*(..., '**PreSilence**') adds further silent frames before the actual start of the sequence.
- *mirevents*(..., '**PostSilence**') adds further silent frames after the actual end of the sequence.
- or the '**Filter**' option: Related options in *mirenvelope* can be specified: '**FilterType**', '**Tau**', '**CutOff**', '**PreDecim**', '**Hilbert**' with same default value than for *mirenvelope*.
  - *mirevents*(..., '**Filterbank**', N) specifies the number of channels for the filterbank decomposition (**mirfilterbank**): the default value being N = 40. N = 0 toggles off the filterbank decomposition.
  - *mirevents*(..., '**FilterbankType**', t) specifies the type of filterbank decomposition (cf. *mirfilterbank*).
  - *mirevents*(..., '**PreSilence**') adds further silence at the beginning of the audio sequence.
  - *mirevents*(..., '**PostSilence**') adds further silence at the end of the audio sequence.
- *mirevents*(..., '**Sum**', 'off') toggles off the channel summation (**mirsum**) that is performed by default.
- Other available options, related to *mirenvelope*: '**HalfwaveCenter**', '**Log**', '**MinLog**', '**Mu**', '**Power**', '**Diff**', '**HalfwaveDiff**', '**Lambda**', '**Center**', '**Smooth**', '**PostDecim**', '**Sampling**', '**UpSample**', all with same default as in *mirenvelope*. In addition, *mirenvelope*'s '**Normal**' option can be controlled as well, with a default set to 1.
- *mirevents*(..., '**SpectralFlux**') computes a spectral flux. Options related to **mirflux** can be passed here as well:
  - '**Inc**' (toggled on by default here),
  - '**Halfwave**' (toggled on by default here),
  - '**Complex**' (toggled off by default as usual),
  - '**Median**' (toggled on by default here, with same default parameters than in *mirflux*).

- *mirevents*(..., '**Emerge**') is an improved version of the '*SpectralFlux*' method that is able to detect more notes and in the same time ignore the spectral variation produced by vibrato.

When the '*Emerge*' method is used for academic research, please cite the following publication:

Lartillot, O., Cereghetti, D., Eliard, K., Trost, W. J., Rappaz, M.-A., Grandjean, D., "Estimating tempo and metrical features by tracking the whole metrical hierarchy", *3rd International Conference on Music & Emotion*, Jyväskylä, 2013.

- *mirevents*(..., '**Pitch**') computes a frame-decomposed autocorrelation function (*mirautocor*), of same default characteristics than those returned by *mirpitch* – with however a range of frequencies set by the following options:

- '**Min**' (set by default to 30 Hz),
- '**Max**' (set by default to 1000 Hz),

and subsequently computes the novelty curve of the resulting similarity matrix. Option related to *mirnovelty* can be passed here as well:

- '**KernelSize**' (set to 32 samples by default).
- *mirevents*(..., '**Novelty**') computes a power-spectrogram (*mirspectrum*) in dB with maximum frequency 1000 Hz, minimal frequency resolution 3 Hz and 50 ms frame with 10 ms hop factor, and subsequently computes the novelty curve of the resulting similarity matrix, using '*Euclidean*' distance and the '*oneminus*' similarity measure. Option related to *mirnovelty* can be passed here as well:
- '**KernelSize**' (set to 64 samples by default).

Besides, the '**Diff**' option can also be used when using the '*Novelty*' method.

Whatever the chosen method, the detection curve is finally converted into an envelope (using *mirenvelope*), and further operations are performed in this order:

- '**Center**' (performed if '*Center*' was specified while calling *mirevents*).
- '**Normal**' (always performed by default).

*mirevents* accepts as input data type either:

- envelope curves (resulting from *mirenvelope*),
- any scalar object, in particular:
  - fluxes (resulting from *mirflux*)
  - novelty curves (resulting from *mirnovelty*)
- similatrix matrices (resulting from *mirsimatrix*): its novelty is automatically computed, with a 'KernelSize' of 32 samples.
- *miraudio* objects, where the audio waveform can be:
  - segmented (using *mirsegment*),
  - decomposed into channels (using *mirfilterbank*),
  - decomposed into frames or not (using *mirframe*):
    - if the audio waveform is decomposed into frames, the detection curve will be based on the spectral flux;
    - if the audio waveform is not decomposed into frames, the default detection curve will be based on the envelope;
- **file name** or the '**Folder**' keyword: same behavior than for *miraudio* objects,
- any other object: it is decomposed into frames (if not already decomposed) using the parameters specified by the 'Frame' option; the flux will be automatically computed by default, or the novelty (if the 'Pitch' option has been chosen).

## EVENT DETECTION

- *miirevents*(..., '**Detect**', *d*) specifies options related to the peak picking from the detection curve:
  - *d* = 'Peaks' (default choice): local maxima are chosen as event positions;
  - *d* = 'Valleys': local minima are chosen as event positions;
  - *d* = 0, or 'no', or 'off': no peak picking is performed.

Options associated to the *mirpeaks* function can be specified as well. In particular:

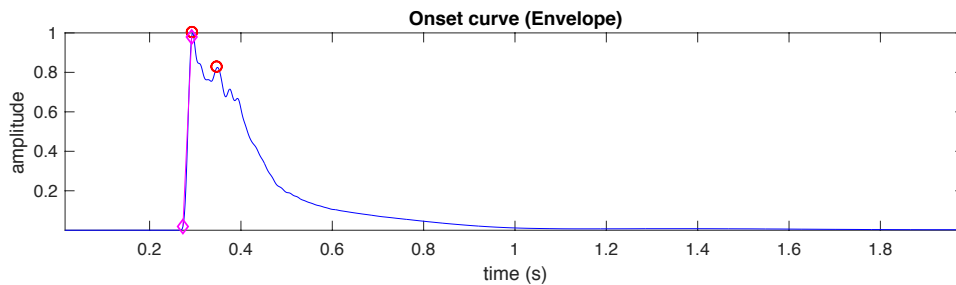
- *miirevents*(..., '**Contrast**', *c*) with default value here  $c = .01$  (but  $c = .05$  if 'Attack' or 'Decay' options are used, cf. next paragraph)
- *miirevents*(..., '**Threshold**', *t*) with default value here  $t = 0$ .


- *mirevents*(..., '**Single**') selects only the highest peak.
- *mirevents*(..., '**SelectFirst**') selects the first out of two closed peaks.
- *mirevents*(..., '**Normalize**', *n*) specifies whether frames are normalized globally or individually. (Same default as in *mirpeaks*).

## ATTACK AND DECAY

The '*Attack*' and '*Decay*' options estimate the beginning and end of the attack and decay phases of each event.

- *mirevents*(..., '**Attack**', *meth*) (or '**Attacks**', *meth*) detects attack phases using the method *meth*, which can be one of the following options: '*Derivate*' (default), '*Effort*' or '*Waveform*'.



*mirevents*('pianoA4.wav', '**Attacks**')  


The '*Derivate*' method works as follows:

- A *slow* detection curve is computed in order to find the events, defined here as the local maxima in the curve. More precisely, for each event we only consider the local minimum preceding the local maximum: it gives a rough estimation of the onset time of the note. This curve is computed using the method chosen by the user ('*Spectro*' (default), '*Filter*', etc.) and with other options set by default (for instance *fl* = .1 s., *fh* = .1 for the '*Spectro*' method) but with '*PreSilence*' and '*PostSilence*' options set as specified by the user (by default, '*PreSilence*' is toggled on while '*PostSilence*' is toggled off).
- A *fast* detection curve is computed in order to find the first local maximum of each note, and to estimate precisely the attack phase. This is this curve that can be controlled by the options when calling *mirevents*, with some changes in the default parameters:
  - If the method '*Spectro*' (default) is chosen, the default frame sizes are this time *fl* = .03 s. and *fh* = .02 (Nymoen et al. 2017).
  - If the chosen method is '*Filter*' instead, the default '*FilterType*' is set to '*Butter*', the '*FilterBank*' option is turned off and the '*Hilbert*' transform is turned on.

- The attack phase is searched for in the temporal region between the onset time (given by the slow curve) and the first local maximum (given by the fast curve). More precisely, in order to reject any low-amplitude peak preceding the attack phase, the search for the attack phase starts at the earliest temporal position where the amplitude of the curve is at 20% of the amplitude of the local maximum.
- The *onset* time, where the attack phase begins, is set at the instant where the rate of the curve (i.e., the value of the first derivate of the curve) is  $o\%$  the maximal rate of the curve in the attack phase, where  $o$  is set by default to 10%, but can be modified using the parameter '**OnsetThreshold**' (expressed as a value between 0 and 1).
- The *attack* time, where the attack phase ends, is set at the instant where the rate of the curve is  $a\%$  the maximal rate of the curve in the attack phase, where  $a$  is set by default to 7.5%, but can be modified using the parameter '**AttackThreshold**' (expressed as a value between 0 and 1).

The '*Effort*' method comes from *Timbre Toolbox* (Peeters et al., 2011).

The parameter '**Alpha**', by default set to 3.75, controls the multiplicative effort ratio. To get the same results as in *Timbre Toolbox*<sup>4</sup>, use the following options:

*mirevents(..., 'Attack', 'Effort', 'Filter', 'CutOff', 5, 'Down', 0, 'Alpha', 3)*

The '*Waveform*' method detects the attack region directly from the audio waveform (Lartillot, Nymoen, Câmara & Danielsen, 2021). The attack region is modeled as a line starting from a low-amplitude point and intersecting one of the local maxima according to two principles: (1) maximizing the slope, while favoring, at the same time, a higher peak if the slope remains only slightly lower and (2) dismissing initial attack regions of relatively low amplitude. The attack start position is fine-tuned by intersecting the attack slope with the audio waveform.

When the '*Waveform*' method is used for academic research, please cite the following publication:

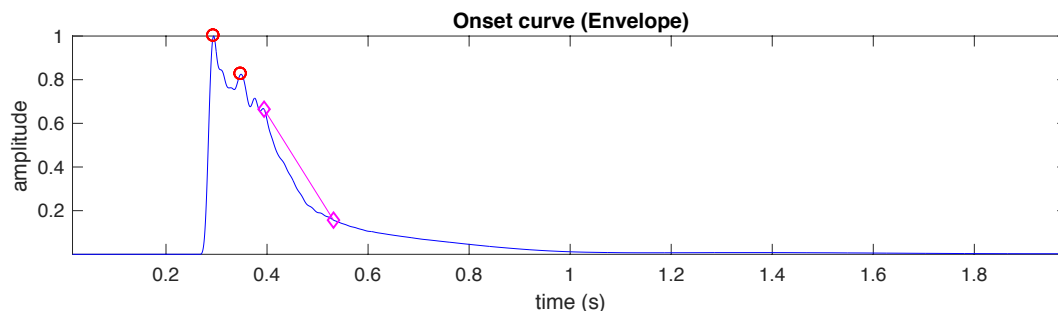
Olivier Lartillot, Kristian Nymoen, Guilherme Schmidt Câmara, & Anne Danielsen (2021). Computational localization of attack regions through a direct observation of the audio waveform, *Journal of the Acoustical Society of America (JASA)*, 149(1), 723—736.

---

<sup>4</sup> One remaining difference is that the *Timbre Toolbox* always normalizes the initial audio files from 0 to 1.



- *mirevents(..., 'Decay')* (or '*Decays*', '*Release*', '*Releases*') detects decay phases, using the '*Derivate*' method, except if the '*Waveform*' option is specified.

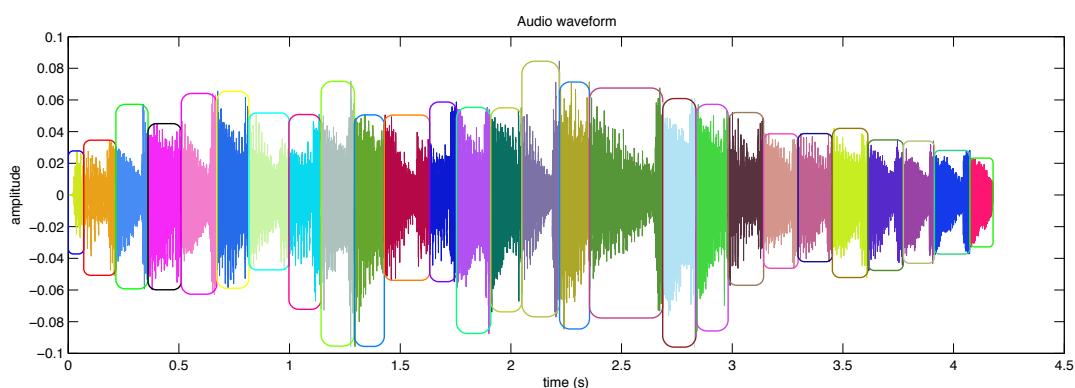


*mirevents('pianoA4.wav', 'Decays')*

- The decay phase is searched for in the temporal region between the last local maximum (given by the fast curve) and the offset time (given by the slow curve). More precisely, in order to reject any low-amplitude peak succeeding the decay phase, the search for the decay phase starts at the latest temporal position where the amplitude of the curve is at 20% of the amplitude of the local maximum.
- The *offset* time, where the decay phase ends, is set at the instant where the rate of the curve (i.e., the value of the first derivate of the curve) is *o*% the maximal rate of the curve in the decay phase, where *o* is set by default to 10%, but can be modified using the parameter '**OffsetThreshold**' (expressed as a value between 0 and 1).
- The *decay* time, where the decay phase starts, is set at the instant where the rate of the curve is *d*% the maximal rate of the curve in the attack phase, where *d* is set by default to 20%, but can be modified using the parameter '**DecayThreshold**' (expressed as a value between 0 and 1).

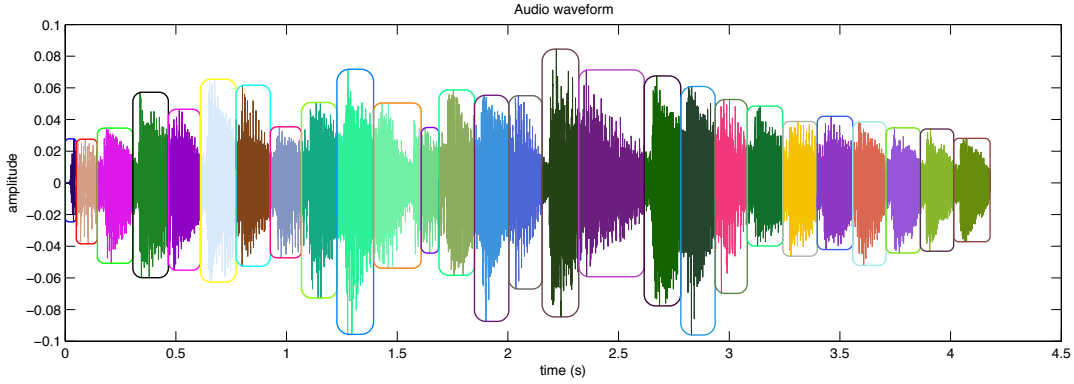
## SEGMENTATION

The temporal localization of events can be used for segmentation of the initial waveform:



```
o = mirevents('ragtime.wav'); mirsegment('ragtime.wav', o)
```

Alternatively, the beginning of the attack phases can be used for the segmentation:



```
o = mirevents('ragtime.wav', 'Attacks'); mirsegment('ragtime.wav', o)
```

## FRAME DECOMPOSITION

The detection curve can be further decomposed into frames if the *'Frame'* option has been specified, with default frame length 3 seconds and hop factor of 10% (0.3 second).

## PRESELECTED MODEL

Complete (or nearly complete) models are available:

- *mirevents(..., 'Scheirer')* follows the model proposed in (Scheirer, 1998). It corresponds to  

```
mirevents(..., 'Filter', 'FilterbankType', 'Scheirer', 'FilterType', 'HalfHann', 'Sampling', 200,  

'HalfwaveDiff', 'Sum', o, 'Detect', o)
```
- *mirevents(..., 'Klapuri99')* follows the model proposed in (Klapuri., 1999). It corresponds to  

```
o = mirevents(..., 'Filter', 'FilterbankType', 'Klapuri', 'FilterType', 'HalfHann', 'PreDecim',  

180, 'Sum', o, 'PostDecim', o);
```

```
o2 = mirenvelope(o, 'HalfwaveDiff'); % absolute distance function D
```

```
o = mirenvelope(o, 'Mu', 'HalfwaveDiff'); % relative distance function W
```

```
p = mirpeaks(o, 'Contrast', .2, 'Chrono');
```

```
p2 = mirpeaks(o2, 'ScanForward', p, 'Chrono');
```

$$o = \text{combinepeaks}(p, p2, .05);$$

where *combinepeaks* is a dedicated function that creates a curve made of burst at position of peaks *p* and with amplitude related to peaks *p2*.

$$o = \text{mirsum}(o, \text{'Weights'}, fB);$$

The intensity is multiplied by the band center frequency *fB*.

$$o = \text{mirenvelope}(o, \text{'Smooth'}, 12);$$

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

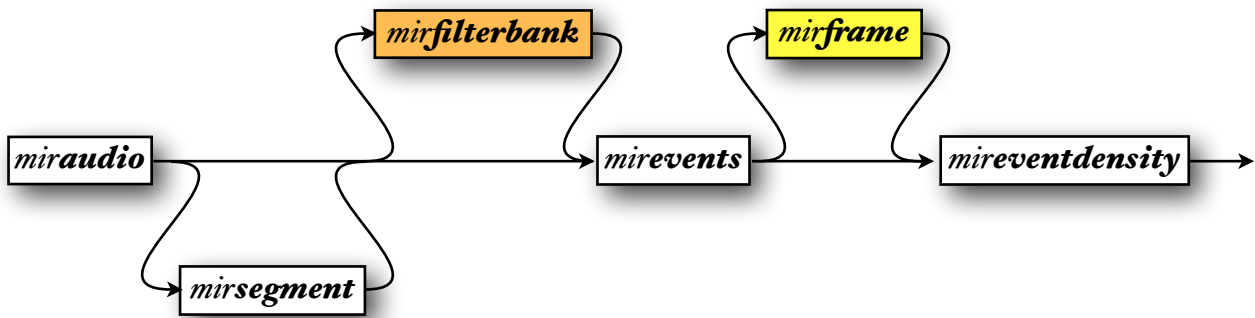
- **'OnsetPos'**: the abscissae position of the starting attack phases, in sample index,
- **'OnsetPosUnit'**: the abscissae position of the starting attack phases, in the default abscissae representation,
- **'AttackPos'**: the abscissae position of the ending attack phases, in sample index,
- **'AttackPosUnit'**: the abscissae position of the ending attack phases, in the default abscissae representation,
- **'DecayPos'**: the abscissae position of the starting decay phases, in sample index,
- **'DecayPosUnit'**: the abscissae position of the starting decay phases, in the default abscissae representation,
- **'OffsetPos'**: the abscissae position of the ending decay phases, in sample index,
- **'OffsetPosUnit'**: the abscissae position of the ending decay phases, in the default abscissae representation.

## *mireventdensity*

### DESCRIPTION

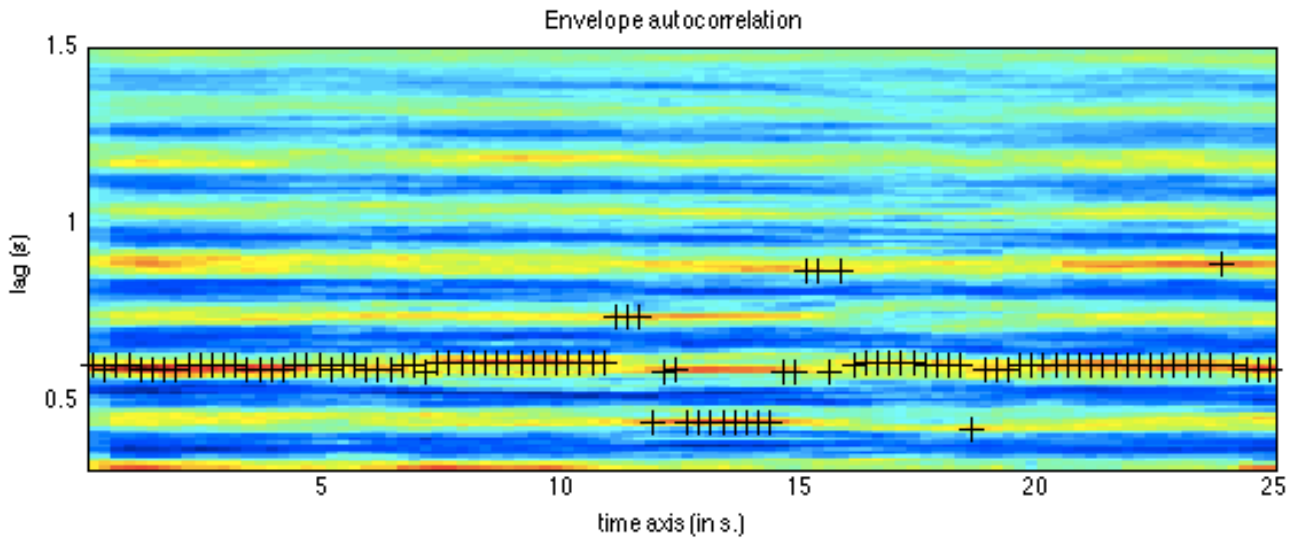
Estimates the average frequency of events, i.e., the number of events detected per second.

### FLOWCHART INTERCONNECTIONS



### FRAME DECOMPOSITION

*mireventdensity*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 10 s and no overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. 'Frame' section.



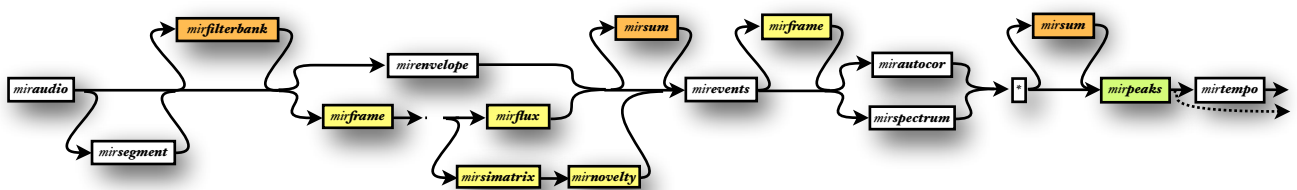
## *mirtempo*

### DESCRIPTION

Estimates the tempo by detecting periodicities from the event detection curve.

#### 1. CLASSICAL VERSION

The classical paradigm for tempo estimation is based on detecting periodicities in a range of BPMs, and choosing the maximum periodicity score for each frame separately.



Tempo estimation is carried out in several steps:

- The event detection curve computed in *mirevents* can be controlled using the following options:
  - **‘Envelope’** (default) and **‘DiffEnvelope’**:
    - with the **‘Method’** set by default to **‘Filter’**:
      - with **‘FilterType’** option with same default,
      - with **‘Filterbank’** option set to 10 by default,
      - with **‘FilterbankType’** option with same default,

- *'Method'* can be set to *'Spectro'* as well, and the *'Freq'*, *'Mel'*, *'Bark'*, *'Cents'* selection can be specified, with same default.
- Besides *'Method'*: *'HalfwaveCenter'*, *'HalfwaveDiff'*, *'Lambda'*, *'Center'*, *'Smooth'*, *'Sampling'*, *'Log'* and *'Mu'*, all with same default and *'Diff'* set to *'On'* by default.
- *'SpectralFlux'*: with *'Complex'*, *'Inc'*, *'Median'* and *'Halfwave'* with same default.
- *'Pitch'* and *'Novelty'*.

Other options related to *mirevents* can be specified:

- *'Filterbank'*, with same default value than for *mirevents*,
- *mirtempo(..., 'Sum',  $w$ )* specifies when to sum the channels. Possible values:
  - $w = \text{'Before'}$ : sum before the autocorrelation or spectrum computation (default).
  - $w = \text{'After'}$ : autocorrelation or spectrum computed for each band, and summed into a "summary".
  - $w = 0$ : tempo estimated for each band separately, with no channel recombination.
- *mirtempo(..., 'Frame', ...)* optionally performs a frame decomposition of the detection curve, with by default a frame length of 3 s and a hop factor of 10% (0.3 s). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.
- Periodicities are detected on the detection curve based on several possible strategies:
  - *mirtempo(..., 'Autocor')* computes an autocorrelation function of the detection curve, using *mirautocor* (default choice). Options related to *mirautocor* can be specified:
    - *'Enhanced'* (toggled on by default<sup>5</sup> here),
    - *'Resonance'* (set by default to *'ToivainenSnyder'*),
    - *'NormalWindow'* (same default value).
  - *mirtempo(..., 'Spectrum')* computes a spectral decomposition of the detection curve, using *mirspectrum*. Options related to *mirspectrum* can be passed here as well:
    - *'ZeroPad'* (set by default here to 10 000 samples),
    - *'Prod'* (same default, when toggled on, as for *mirspectrum*),
    - *'Resonance'* either *'ToivainenSnyder'* (default value) or 0, *'off'*, or *'no'*.

---

<sup>5</sup> except when *'Track'* option is used, as explained.

- *mirtempo*(..., '**Autocor**', '**Spectrum**') combines both strategies: the autocorrelation function is translated into the frequency domain in order to be compared to the spectrum curve, and the two curves are subsequently multiplied.
- A peak picking is applied to the autocorrelation function or to the spectrum representation. The parameters of the peak picking can be tuned:
  - *mirtempo*(..., '**Total**', *m*) selects not only the best tempo, but the *m* best tempos.
  - *mirtempo*(..., '**Min**', *mi*) indicates the lowest tempo taken into consideration, expressed in bpm. Default value: 40 bpm.
  - *mirtempo*(..., '**Max**', *ma*) indicates the highest tempo taken into consideration, expressed in bpm. Default value: 200 bpm.
  - *mirtempo*(..., '**Track**', *t*) tracks peaks along time, in order to obtain a stabilized tempo curve and to limit therefore switches between alternative pulsations. Default value when option toggled on:  $t = 0.1$  s. When 'Track' is toggled on, 'Enhanced' is forced to off.

When 'Track' is used for academic research, please cite the following publication:

Olivier Lartillot, "[mirtempo: Tempo estimation through advanced frame-by-frame peaks tracking](#)", *Music Information Retrieval Evaluation eXchange (MIREX 2010)*.

- *mirtempo*(..., '**Contrast**', *c*) specifies the contrast factor for the peak picking. Default value:  $c = 0.1$
- *mirtempo*(..., '**Nearest**', *n*) chooses the peak closest to *n* (in s.). Default value when option toggled on:  $n = 0.5$  s.

*mirtempo* accepts as input data type either:

- *mirautocor* objects,
- *mirspectrum* objects,
- event detection curve (resulting from *mirevents*), frame-decomposed or not, channel-decomposed or not,
- and all the input data accepted by *mirevents*.

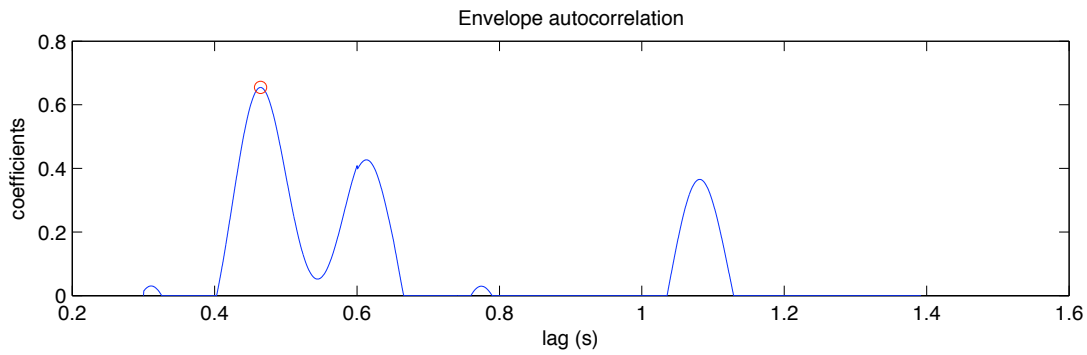
*mirtempo* can return several outputs:

1. the tempo itself (or set of tempi) and
2. the ***mirspectrum*** or ***mirautocor*** data, where is highlighted the (set of) peak(s) corresponding to the estimated tempo (or set of tempi).

The tempo estimation related to the *ragtime.wav* example

$$[t \ ac] = \text{mirtempo}(\text{'ragtime.wav'})$$

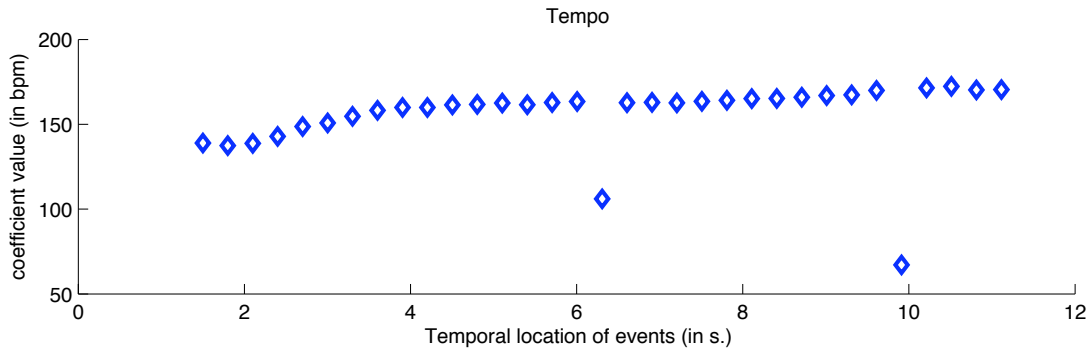
leads to a tempo  $t = 129.1832$  bpm and to the following autocorrelation curve *ac*:



The frame-decomposed tempo estimation related to the *czardas* example

$$[t \ ac] = \text{mirtempo}(\text{'czardas'}, \text{'Frame'})$$

leads to the following tempo curve *t*:

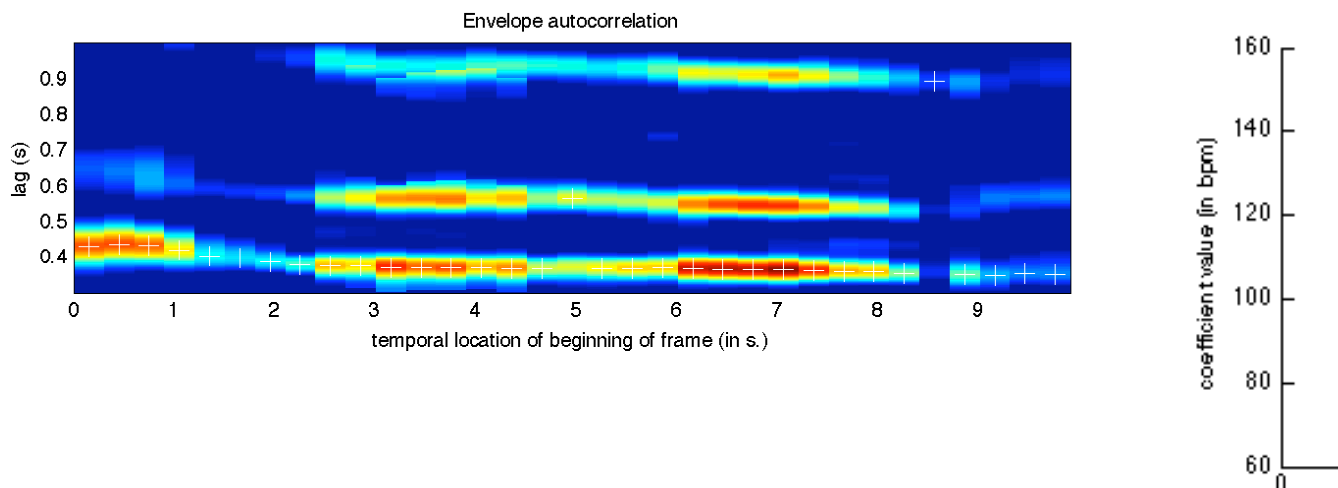


and the following autocorrelation frame decomposition *ac*:

Below are the results of the analysis of a more challenging example: the first seconds of the first movement of a performance of J.S. Bach's *Brandenburg concert No.2 in F Major*, BWV 1047:

The classical method generates a tempo curve with a lot of shifts from one metrical level to another.





## 2 . M E T R E - B A S E D V E R S I O N

*mirtempo*(..., '**Metre**') tracks tempo by building a hierarchical metrical structure (using *mirmetre*). This enables to find coherent metrical levels leading to a continuous tempo curve.

When the '*Metre*' option is used for academic research, please cite the following publication:

Lartillot, Olivier & Grandjean, Didier (2019). Tempo and Metrical Analysis by Tracking Multiple Metrical Levels Using Autocorrelation. *Applied Sciences*, 9(23). doi: 10.3390/app9235121.

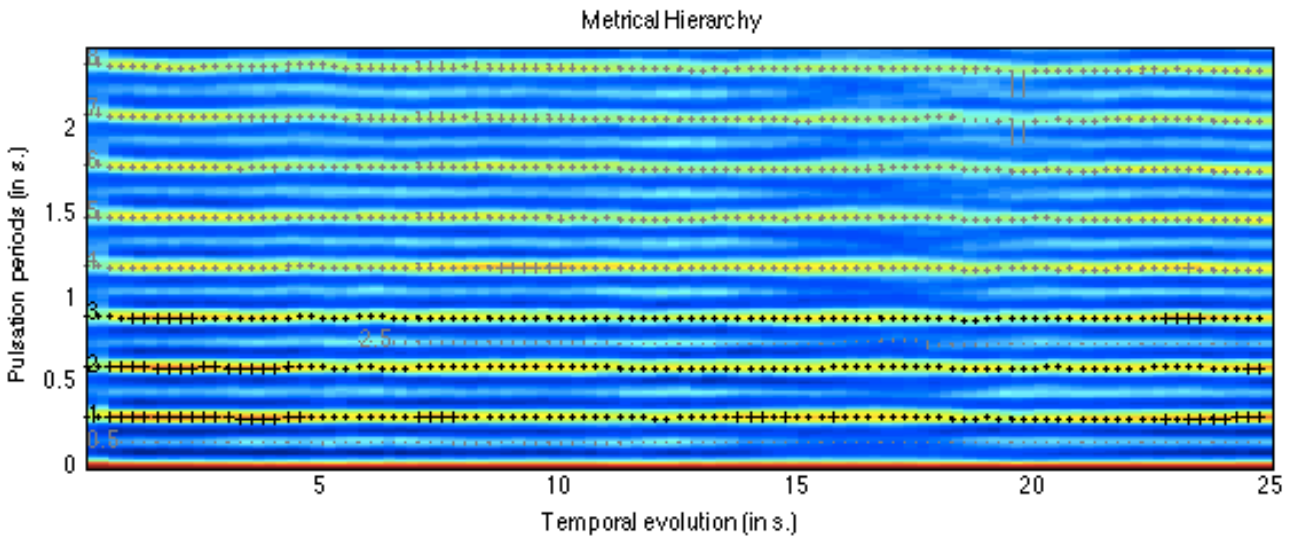
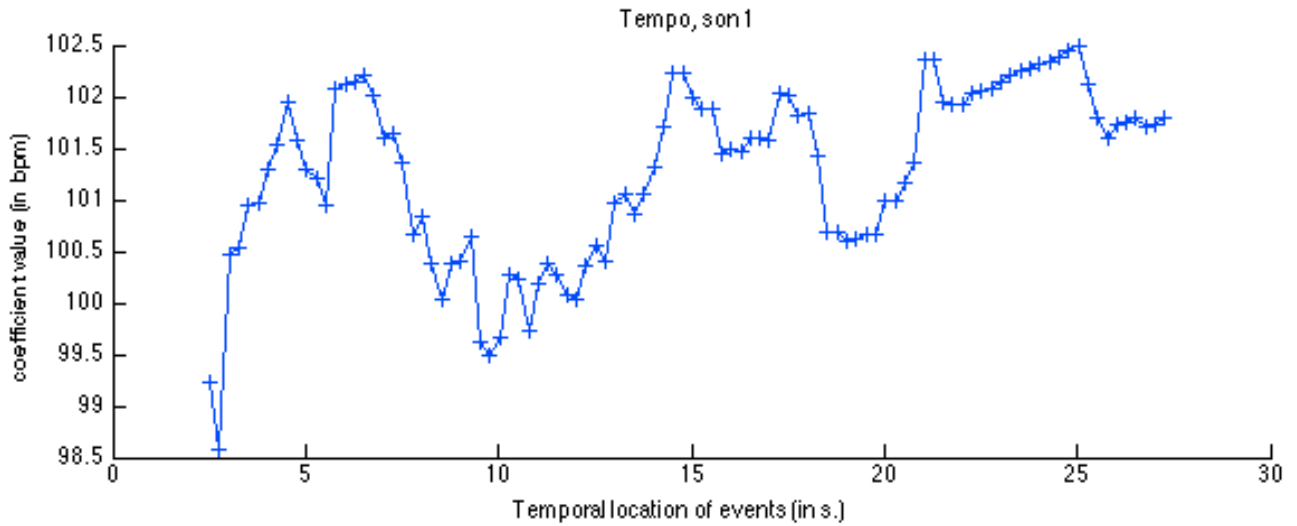
*mirtempo*(..., '*Metre*') accepts as input data type either:

- *mirautocor* objects,
- event detection curve (resulting from *mirevents*), frame-decomposed or not, channel-decomposed or not,
- and all the input data accepted by *mirevents*.

*mirtempo*(..., '*Metre*') can return several outputs:

1. the tempo itself and
2. the ***mirmetre*** representation.

Below are the results of the analysis of the same excerpt of J.S. Bach's *Brandenburg concert No.2 in F Major*, BWV 1047, this time using the '*Metre*' option:



The metrical structure built using the ‘*Metre*’ strategy enables to find coherent metrical levels leading to a continuous tempo curve.

## TEMPO CHANGE

*mirtempo*(..., ‘**Change**’) computes the difference between successive values of the tempo curve. Tempo change is expressed independently from the choice of a metrical level by computing the ratio of tempo values between successive frames, and is expressed in logarithmic scale (base 2), so that no tempo change gives a value of 0, increase of tempo gives positive value, and decrease of tempo gives negative value.

## *mirmetre*

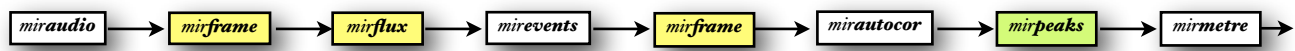
### METRICAL ANALYSIS

Provides a detailed description of the hierarchical metrical structure by detecting periodicities from the event detection curve and tracking a broad set of metrical levels.

When *mirmetre* is used for academic research, please cite the following publication:

Lartillot, Olivier & Grandjean, Didier (2019). Tempo and Metrical Analysis by Tracking Multiple Metrical Levels Using Autocorrelation. *Applied Sciences*, 9(23). doi: 10.3390/app9235121.

### FLOWCHART INTERCONNECTIONS



The metrical hierarchy is estimated through the following steps:

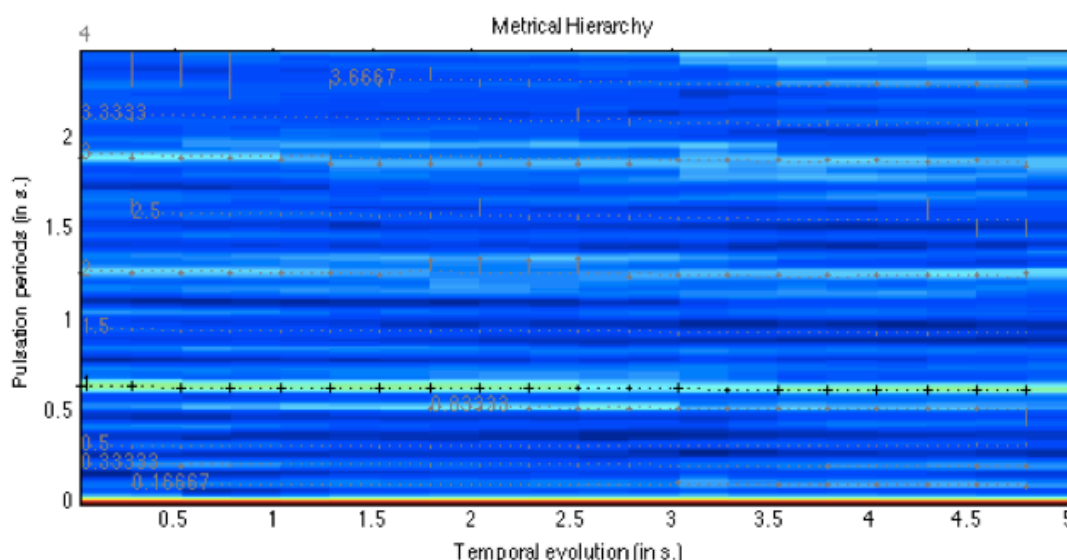
- An event detection curve is computed using *mirevents* with the ‘*Emerge*’ method and no actual event detection (‘*Detect*’ set to 0).
- An autocorrelation function is computed, using *mirautocor*, on the successive frames of the event detection curve. Default options are used apart from those specified below:
  - *mirmetre*(..., ‘**Frame**’, ...) specifies the parameters of the frame decomposition, with by default a frame length of 5 s and a hop factor of 5% (0.25 s).
  - *mirmetre*(..., ‘**Min**’, *mi*) indicates the lowest periodicity taken into consideration, expressed in bpm. Default value: 24 bpm.
  - *mirmetre*(..., ‘**Max**’, *ma*) specifies the highest periodicity taken into consideration, expressed in bpm. Default value: 1000.
- A peak picking is applied to the autocorrelation function, taking all peaks satisfying the constraints specified by the options below:
  - *mirmetre*(..., ‘**Contrast**’, *c*) specifies the contrast factor for the peak picking. Default value:  $c = 0.05$ .
  - *mirmetre*(..., ‘**Threshold**’, *c*) specifies the contrast factor for the peak picking. Default value:  $c = 0$ .

*mirmetre* accepts as input data type either:

- *mirmetre* objects,
- *mirautocor* objects,
- event detection curve (resulting from *mirevents*), frame-decomposed or not, channel-decomposed or not,
- and all the input data accepted by *mirevents*.

## EXAMPLE

The figure below shows the output of the metrical analysis of the first five seconds of a performance of the *Finale (Allegro energico)* of M. Bruch's *Violin Concerto No.1 in G minor, op.26*. On top of the autocorrelogram (corresponding to a standard output of the frame-decomposed *mirautocor* computation), the metrical structure as tracked by the algorithm is annotated.



Metrical levels are indicated with lines of crosses that are linked together between successive frames with dotted lines. The level index is indicated on the left of each line. The dominant metrical level, indicated as level 1, is drawn in black, while other levels are shown in light brown. The cross size indicates the related pulse strength, corresponding to the autocorrelation score for that periodicity. If the actual periodicity is deviated from the theoretical harmonic series of periodicities expected from a metrical structure, a vertical line is drawn between the actual and the theoretical periods.

In this example, level 1 is subdivided into six sub-beats, with its elementary level 1/6, as well as its half slower level 2/6 corresponding to a ternary division of level 1, and finally its three times slower level 3/6 corresponding to a binary division of level 1.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- ***Autocor***: the autocorrelogram,
- ***Globpm***: the series of BPM values for each successive frame, related to level 1 of each metrical structure.

## *mirmetroid*

### METRICAL CENTROID AND STRENGTH

Provides two description derived from the metrical analysis carried out using *mirmetre*, with a further selection of primary metrical levels.

1. *Dynamic metrical centroid*: an assessment of metrical activity that is based on the computation of the centroid of the selected metrical levels. The resulting metrical centroid curve indicates the temporal evolution of the metrical activity expressed in BPM, so that the values can be compared with the tempo values also in BPM. High BPM values for the metrical centroid indicate that more elementary metrical levels (i.e., very fast levels corresponding to very fast rhythmical values) predominate. Low BPM values indicate on the contrary that higher metrical levels (i.e., slow pulsations corresponding to whole notes, bars, etc.) predominate. If one particular level is particularly dominant, the value of the metrical centroid naturally approaches the corresponding tempo value on that particular level.
2. *Dynamic metrical strength*: indicates whether there is a clear and strong pulsation, or even a strong metrical hierarchy, or whether on the other hand the pulsation is somewhat hidden, unclear, or there is a complex mixture of pulsations. Instead of simply identifying beat strength with the autocorrelation score of that main metrical level, we sum the autocorrelation scores of the selected metrical levels. The metrical strength is increased by any increase of autocorrelation score at any dominant level, or if new dominant levels are added to the selection. Whereas the autocorrelation score is a value lower than 1, the metrical strength can exceed 1.

When *mirmetroid* is used for academic research, please cite the following publication:

Lartillot, Olivier & Grandjean, Didier (2019). Tempo and Metrical Analysis by Tracking Multiple Metrical Levels Using Autocorrelation. *Applied Sciences*, 9(23). doi: 10.3390/app9235121.

### FLOWCHART INTERCONNECTIONS

*mirmetroid* accepts as input data type either:

- *mirmetre* objects,
- *mirautocor* objects,
- event detection curve (resulting from *mirevents*), frame-decomposed or not, channel-decomposed or not,
- and all the input data accepted by *mirevents*.

## OPTIONS

For each frame, the dominant metrical levels are selected and the centroid of their periodicity (in seconds) is computed, using as weights the amount of autocorrelation score at that specific level that exceeds the autocorrelation score of the underlying metrical sub-levels. In this way, any sudden change in the number of selected metrical levels from one time frame to the successive one does not lead to abrupt changes in the metrical centroid curve.

- *mirmetroid*(..., '**Gate**') uses a simpler method, where the weights are simply identified to the corresponding autocorrelation scores, leading to possible abrupt changes in the metrical centroid curve.

If several metrical hierarchies are detected in parallel, separate metrical centroid and strength curves for each individual metrical hierarchy. The multiple metrical centroid and strength curves are then combined into one metrical centroid curve and one metrical strength curve. The combined metrical centroid curve is the weighted average of the metrical centroid curves of the separate metrical hierarchies, using the metrical strengths as weights. The combined metrical strength is the summation of the metrical strength curves of the separate metrical hierarchies.

- *mirmetroid*(..., '**Combine**', *o*) does not perform the combination of the metrical centroid and strength curves, displaying therefore as many curves as metrical hierarchies detected.

## mirpulseclarity

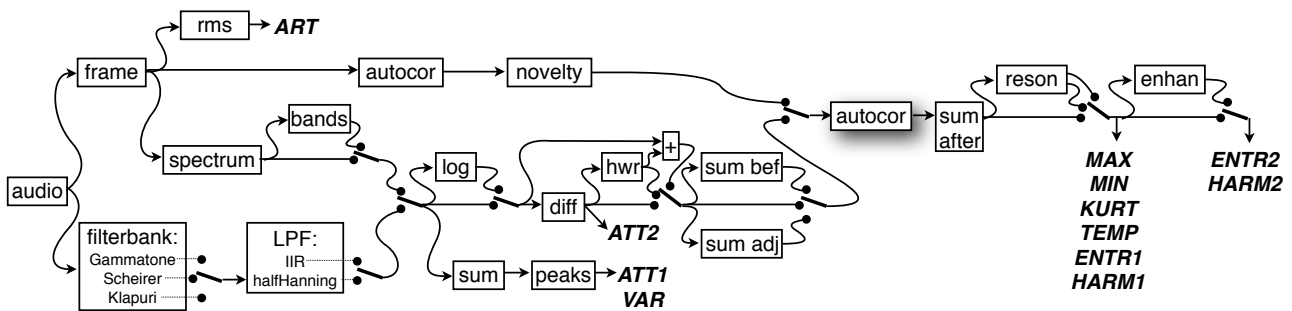
When *mirpulseclarity* is used for academic research, please cite the following publication:

Olivier Lartillot, Tuomas Eerola, Petri Toiviainen, Jose Fornari, "Multi-feature modeling of pulse clarity: Design, validation, and optimization", *International Conference on Music Information Retrieval*, Philadelphia, 2008.

### DESCRIPTION

Estimates the rhythmic clarity, indicating the strength of the beats estimated by the *mirtempo* function.

### FLOWCHART INTERCONNECTIONS



The pulse clarity can be estimated in various ways:

- *mirpulseclarity*(..., *s*) selects a particular heuristic for pulse clarity estimation. Most heuristics are based on the autocorrelation curve computed for tempo estimation (i.e., the second output of *mirtempo*) (Lartillot, Eerola, Toiviainen, and Fornari, 2008):
  - *s* = '**MaxAutocor**' selects the maximum correlation value in the autocorrelation curve (default heuristic).
  - *s* = '**MinAutocor**' selects the minimum correlation value in the autocorrelation curve.
  - *s* = '**MeanPeaksAutocor**' averages the local maxima in the autocorrelation curve.
  - *s* = '**KurtosisAutocor**' computes the kurtosis of the autocorrelation curve.
  - *s* = '**EntropyAutocor**' computes the entropy of the autocorrelation curve.
  - *s* = '**InterfAutocor**' computes the harmonic relations between pulsations.
  - *s* = '**TempoAutocor**' selects the tempo related to the highest autocorrelation.

Others heuristics are based more simply on the detection curve itself:



- $s = \text{'Articulation'}$  estimates the average silence ratio of the detection curve (option *'ASR'* in *mirlowenergy*).
- $s = \text{'Attack'}$  averages the attack slopes of all events (the *'Diff'*, *'Gauss'* can be specified, with same default).
- $s = \text{'ExtremEnvelope'}$  estimates the total amplitude variability of the detection curve.

*mirpulseclarity(..., 'Model', m)* selects one out of two possible models that have been found as optimal in our experiments (Lartillot, Eerola, Toivianen, and Fornari, 2008):

- $m = 1$  selects the default model with its associated weight.
- $m = 2$  selects the following model: *'Gammatone'*, no log, no *'Resonance'*, *'Lambda'* set to .8, and *'Sum'* set to *'After'*, with its associated weight.
- $m = [1\ 2]$  sums the two models altogether.

The event detection curve computed in *miरेvents* can be controlled using the following options:

- *'Envelope'* (default) and *'DiffEnvelope'*:
  - with the *'Method'* set by default to *'Spectro'*, and the *'Freq'*, *'Mel'*, *'Bark'*, *'Cents'* selection can be specified, with same default.
  - *'Method'* can be set to *'Filter'* as well:
    - with *'FilterType'* option with same default,
    - with *'Filterbank'* option set to 20 by default,
    - with *'FilterbankType'* option set to *'Scheirer'* by default,
  - Besides *'Method'*: *'HalfwaveDiff'*, *'Lambda'*, *'Smooth'*, *'Log'* with same default, and *'Mu'*, set by default here to 100.
- *'SpectralFlux'*: with *'Inc'* with same default, and *'Median'* and *'Halfwave'* toggled off by default.
- and *'Pitch'*.

The autocorrelation function performed in *mirautocor* can be controlled using the following options:

- *'Enhanced'* (toggled off by default; forced to *'Off'* in *'MinAutocor'*),
- *'Resonance'*, *'Min'*, *'Max'* (with same default as in *mirautocor*).

Some further options operates as in *mirtempo*:

- **'Sum'**,
- **'Total'** (ignored in *'MaxAutocor'*, *'MinAutocor'* and *'EntropyAutocor'* methods),
- **'Contrast'**: with a default value set to .01,

*mirpulseclarity* accepts as input data type either:

- *mirautocor* objects,
- event detection curve (resulting from *mirevents*), frame-decomposed or not, channel-decomposed or not,
- and all the input data accepted by *mirevents*.

*mirpulseclarity* can return several outputs:

1. the pulse clarity value and
2. the *mirautocor* data that was used for the estimation of pulse clarity.

## FRAME DECOMPOSITION

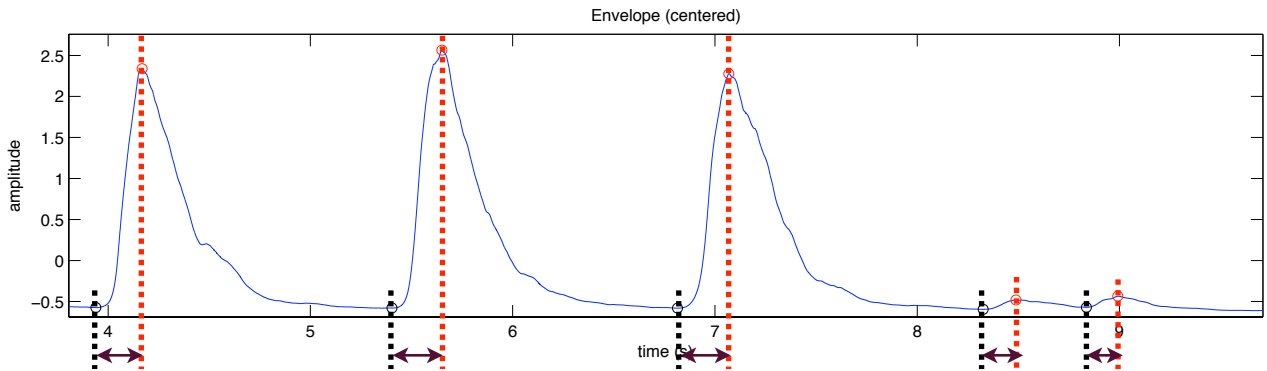
*mirpulseclarity*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 5 s and a hop factor of 10% (0.5 s). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

### 3.3. Timbre

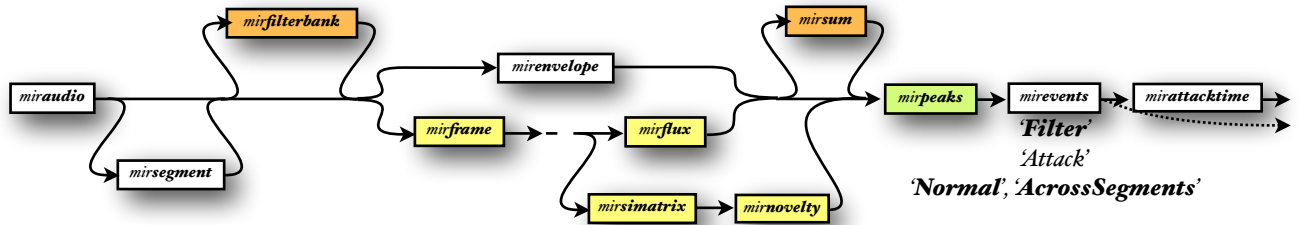
## *mirattacktime*

### DESCRIPTION

The attack phase detected using the ‘Attacks’ option in *mirevents* can offer some timbral characterizations. One simple way of describing the attack phase, proposed in *mirattacktime*, consists in estimating its temporal duration.



### FLOWCHART INTERCONNECTIONS



*mirattacktime* accepts as input data type either:

- event detection curves (resulting from *mirevents*), already including peaks or not,
- and all the input data accepted by *mirevents*.

*mirevents* is normalized using the ‘**Normal**’ option set to ‘**AcrossSegments**’.

Some options in *mirevents* can be controlled:

- *mirattacktime*(..., ‘**Filter**’) uses the ‘Filter’ method instead of the default ‘Spectro’.
- *mirattacktime*(..., ‘**Down**’, *r*) controls the ‘Down’ option in *mirevents*.
- *mirattacktime*(..., ‘**CutOff**’, *f*) controls the ‘CutOff’ option in *mirevents*.
- *mirattacktime*(..., ‘**Single**’) toggles on the ‘Single’ option in *mirevents*.

- *mirattacktime*(..., '**LogCurve**') toggles on the 'Log' option in *mirevents*.
  - *mirattacktime*(..., '**MinLog**', *ml*) controls the 'MinLog' option in *mirevents*.
- *mirattacktime*(..., '**Contrast**', *c*) controls the 'Contrast' option in *mirevents*.
- *mirattacktime*(..., '**PreSilence**', *o*) toggles off the 'PreSilence' option in *mirevents*.
- *mirattacktime*(..., '**PostSilence**') toggles on the 'PostSilence' option in *mirevents*.
- *mirattacktime*(..., '**Attacks**', *meth*) controls the 'Attacks' option in *mirevents*.

*mirattacktime* can return several outputs:

1. the attack time itself and
2. the event detection curve returned by *mirevents*, including the detected events and the attack phases.

## OPTIONS

- *mirattacktime*(..., *scale*) specifies the output scale, linear or logarithmic. Possible values for *scale* are:
  - *scale* = '**Lin**' returns the duration in a linear scale (in seconds). (Default choice)
  - *scale* = '**Log**' returns the duration in a log scale (Krimphoff et al., 1994).

The log attack time given by the *Timbre Toolbox* (Peeters et al., 2011) can be obtained by using the following set of options:

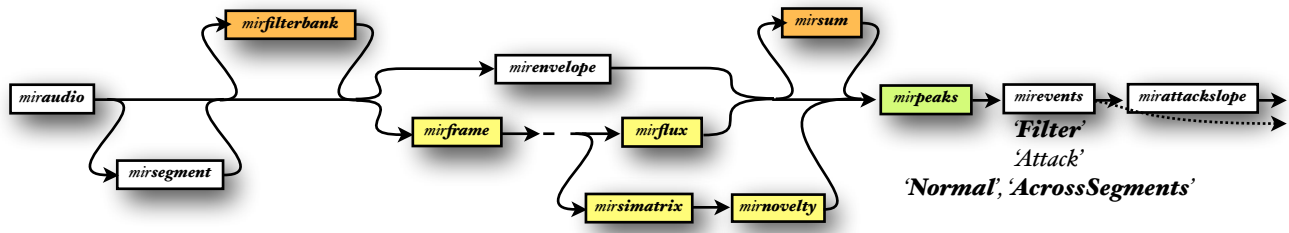
*mirattacktime*(..., '*Filter*', '*Attack*', '*Effort*', '*Log*', '*Down*', 0, '*CutOff*', 5)

## *mirattackslope*

### DESCRIPTION

Another description of the attack phase is related to its average slope. Values are expressed in the same scale than the original signal, but normalized in time (expressed in seconds).

### FLOWCHART INTERCONNECTIONS



*mirattackslope* accepts as input data type either:

- event detection curves (resulting from *mirevents*),
- and all the input data accepted by *mirevents*.

Some options in *mirevents* can be controlled:

- *mirattackslope*(..., **Filter**) uses the 'Filter' method instead of the default 'Spectro'.
  - *mirattackslope*(..., **Down**, *r*) controls the 'Down' option in *mirevents*.
- *mirattackslope*(..., **Single**) toggles on the 'Single' option in *mirevents*.
- *mirattackslope*(..., **LogCurve**) toggles on the 'Log' option in *mirevents*.
  - *mirattackslope*(..., **MinLog**, *ml*) controls the 'MinLog' option in *mirevents*.
- *mirattackslope*(..., **Normal**, *n*) controls the 'Normal' option in *mirevents* (by default set to 'AcrossSegments').
- *mirattackslope*(..., **PreSilence**, *o*) toggles off the 'PreSilence' option in *mirevents*.
- *mirattackslope*(..., **PostSilence**) toggles on the 'PostSilence' option in *mirevents*.
- *mirattackslope*(..., **Attacks**, *meth*) controls the 'Attacks' option in *mirevents*.

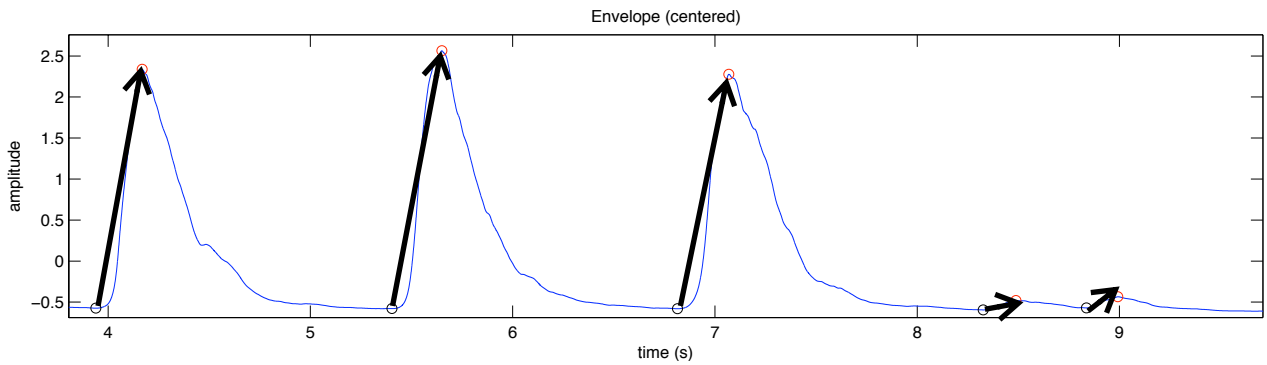
The peak picking from the event detection is performed in any case. Its **Contrast** parameter can be specified. Its default value is the same as in *mirevents*.

*mirattackslope* can return several outputs:

1. the attack slope itself and
2. the event detection curve returned by *mirevents*, including the detected events and the attack phases.

## OPTIONS

- *mirattackslope*(*x*, *meth*) specifies the method for slope estimation. Possible values for *meth* are:
  - *meth* = '**Diff**' computes the slope as a ratio between the magnitude difference at the beginning and the ending of the attack period, and the corresponding time difference. (Default choice)

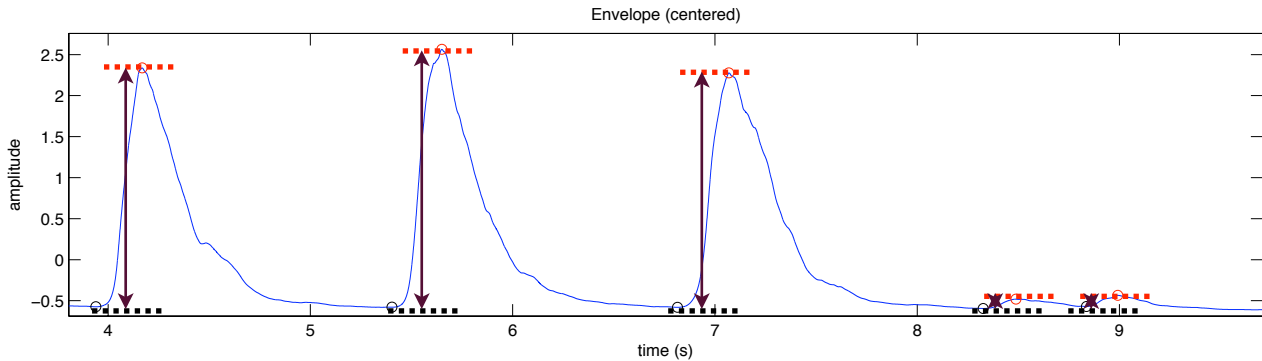


- *meth* = '**Gauss**' computes the average of the slope, weighted by a gaussian curve that emphasizes values at the middle of the attack period (similar to Peeters, 2004).

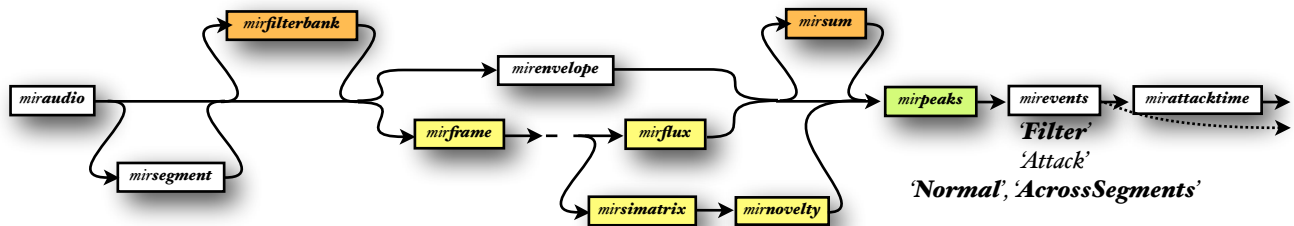
## mirattackleap

### DESCRIPTION

Another simple way of describing the attack phase, proposed in *mirattackleap*, consists in estimating the amplitude difference between the beginning and the end of the attack phase. Values are expressed in the same scale than the original signal.



### FLOWCHART INTERCONNECTIONS



*mirattackleap* accepts as input data type either:

- event detection curves (resulting from *mirevents*), already including peaks or not,
- and all the input data accepted by *mirevents*.

Some options in *mirevents* can be controlled:

- *mirattackleap*(..., **Filter**) uses the 'Filter' method instead of the default 'Spectro'.
  - *mirattackleap*(..., **Down**, *r*) controls the 'Down' option in *mirevents*.
- *mirattackleap*(..., **Single**) toggles on the 'Single' option in *mirevents*.
- *mirattackleap*(..., **LogCurve**) toggles on the 'Log' option in *mirevents*.
  - *mirattackleap*(..., **MinLog**, *ml*) controls the 'MinLog' option in *mirevents*.
- *mirattackleap*(..., **Normal**, *n*) controls the 'Normal' option in *mirevents* (by default set to 'AcrossSegments').

- *mirattackleap*(..., '**PreSilence**', o) toggles off the '*PreSilence*' option in *mirevents*.
- *mirattackleap*(..., '**PostSilence**') toggles on the '*PostSilence*' option in *mirevents*.
- *mirattackleap*(..., '**Attacks**', *meth*) controls the '*Attacks*' option in *mirevents*.

The peak picking from the event detection is performed in any case. Its '**Contrast**' parameter can be specified. Its default value is the same as in *mirevents*.

*mirattackleap* can return several outputs:

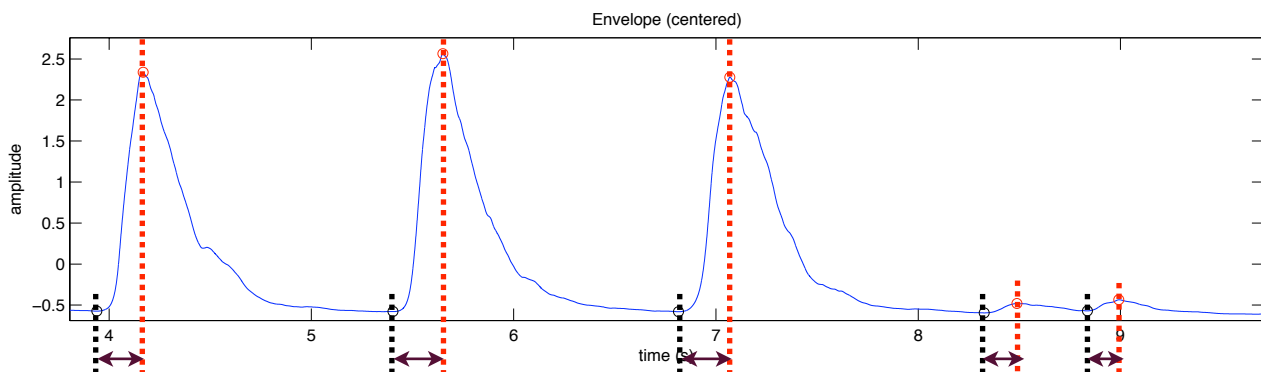
1. the attack leap itself and
2. the event detection curve returned by *mirevents*, including the detected events and the attack phases.



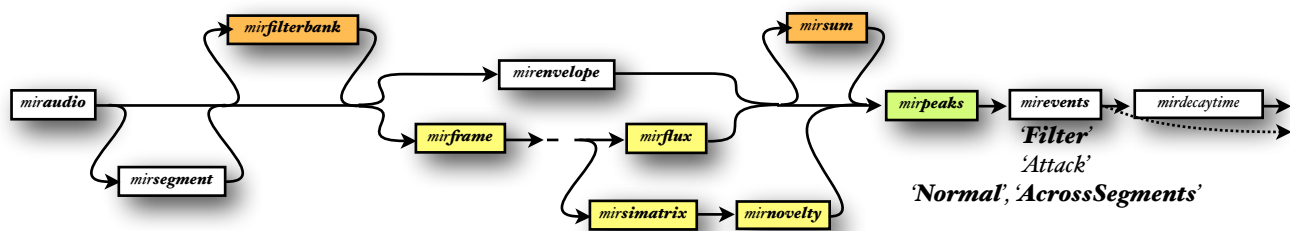
## mirdecaytime

### DESCRIPTION

Description of the decay phase related to its temporal duration, similar to *mirattacktime*.



### FLOWCHART INTERCONNECTIONS



*mirdecaytime* accepts as input data type either:

- event detection curves (resulting from *mirevents*), already including peaks or not,
- and all the input data accepted by *mirevents*. In this case *mirevents* is called with the '**Spectro**' method.

*mirevents* is normalized using the '**Normal**' option set to '**AcrossSegments**'.

Some options in *mirevents* can be controlled:

- *mirdecaytime*(..., '**Single**') toggles on the 'Single' option in *mirevents*.
- *mirdecaytime*(..., '**LogCurve**') toggles on the 'Log' option in *mirevents*.
- *mirdecaytime*(..., '**MinLog**', *ml*) controls the 'MinLog' option in *mirevents*.
- *mirdecaytime*(..., '**Contrast**', *c*) controls the 'Contrast' option in *mirevents*.
- *mirdecaytime*(..., '**PreSilence**', *o*) toggles off the 'PreSilence' option in *mirevents*.
- *mirdecaytime*(..., '**PostSilence**') toggles on the 'PostSilence' option in *mirevents*.

*mirdecaytime* can return several outputs:

1. the decay time itself and
2. the event detection curve returned by *mirevents*, including the detected events and the decay phases.

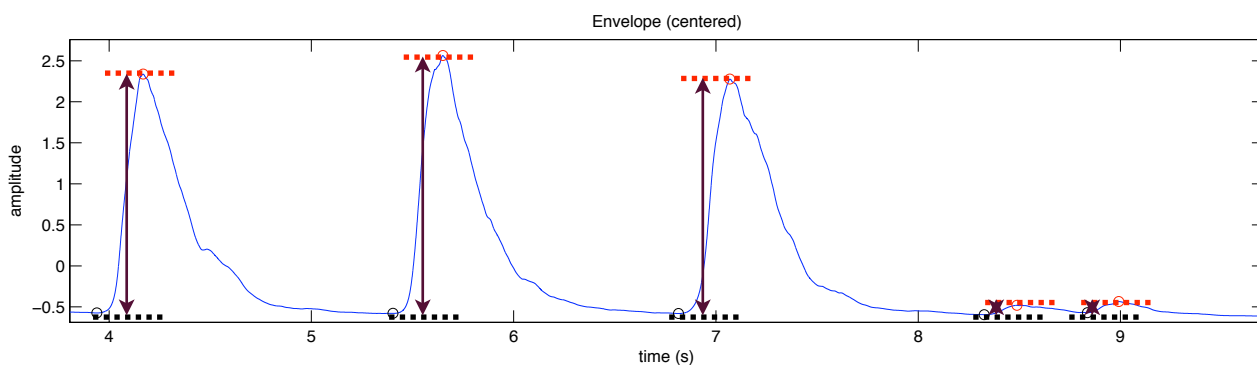
## OPTIONS

- *mirdecaytime*(..., *scale*) specifies the output scale, linear or logarithmic. Possible values for *scale* are:
  - *scale* = '**Lin**' returns the duration in a linear scale (in seconds). (Default choice)
  - *scale* = '**Log**' returns the duration in a log scale.

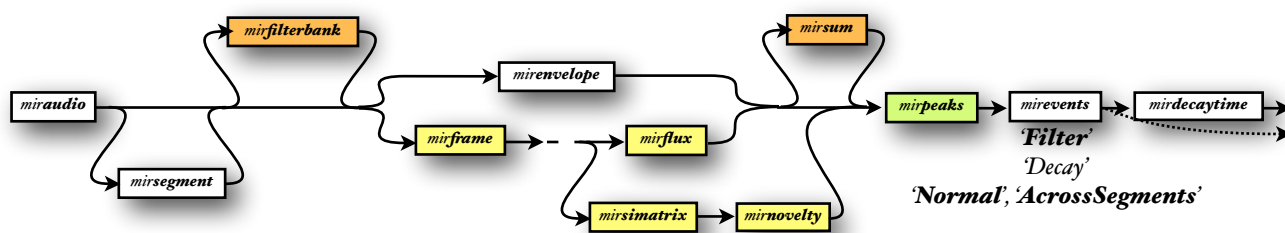
## mirdecayleap

### DESCRIPTION

Description of the decay phase related to its amplitude change, similar to *mirattackleap*. Values are expressed in the same scale than the original signal.



### FLOWCHART INTERCONNECTIONS



*mirdecayleap* accepts as input data type either:

- event detection curves (resulting from *mirevents*), already including peaks or not,
- and all the input data accepted by *mirevents*. In this case *mirevents* is called with the '**Spectro**' method.

Some options in *mirevents* can be controlled:

- *mirdecayleap*(..., '**Single**') toggles on the 'Single' option in *mirevents*.
- *mirdecayleap*(..., '**LogCurve**') toggles on the 'Log' option in *mirevents*.
- *mirdecayleap*(..., '**MinLog**', *ml*) controls the 'MinLog' option in *mirevents*.
- *mirdecayleap*(..., '**Normal**', *n*) controls the 'Normal' option in *mirevents* (by default set to 'Across-Segments').
- *mirdecayleap*(..., '**PreSilence**', *o*) toggles off the 'PreSilence' option in *mirevents*.
- *mirdecayleap*(..., '**PostSilence**') toggles on the 'PostSilence' option in *mirevents*.

The peak picking from the event detection is performed in any case. Its '**Contrast**' parameter can be specified. Its default value is the same as in *mirevents*.

*mirdecayleap* can return several outputs:

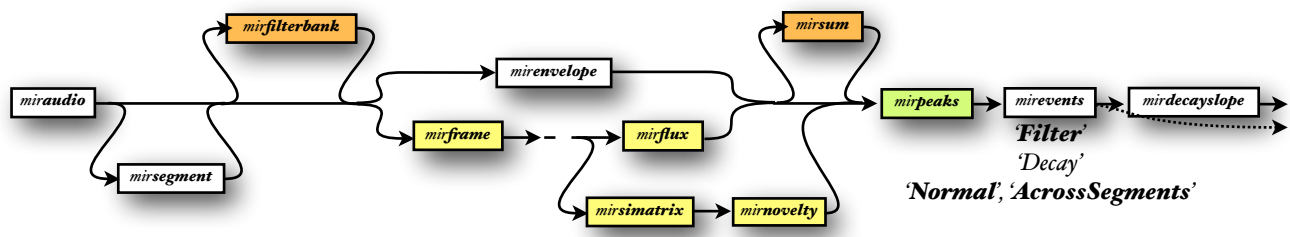
1. the decay leap itself and
2. the event detection curve returned by *mirevents*, including the detected events and the decay phases.

## *mirdecayslope* (previously *mirdecreaseslope*)

### DESCRIPTION

Description of the decay phase related to its decay slope, similar to *mirattackslope*. Values are expressed in the same scale than the original signal, but normalized in time (expressed in seconds).

### FLOWCHART INTERCONNECTIONS



*mirdecayslope* accepts as input data type either:

- event detection curves (resulting from *mirevents*),
- and all the input data accepted by *mirevents*. In this case *mirevents* is called with the '**Spectro**' method.

Some options in *mirevents* can be controlled:

- *mirdecayslope*(..., '**Single**') toggles on the 'Single' option in *mirevents*.
- *mirdecayslope*(..., '**LogCurve**') toggles on the 'Log' option in *mirevents*.
- *mirdecayslope*(..., '**MinLog**', *ml*) controls the 'MinLog' option in *mirevents*.
- *mirdecayslope*(..., '**Normal**', *n*) controls the 'Normal' option in *mirevents* (by default set to 'AcrossSegments').
- *mirdecayslope*(..., '**PreSilence**', *o*) toggles off the 'PreSilence' option in *mirevents*.
- *mirdecayslope*(..., '**PostSilence**') toggles on the 'PostSilence' option in *mirevents*.

The peak picking from the event detection is performed in any case. Its '**Contrast**' parameter can be specified. Its default value is the same as in *mirevents*.

*mirdecayslope* can return several outputs:

1. the decay slope itself and

2. the event detection curve returned by ***mirevents***, including the detected events and the decay phases.

## OPTIONS

- *mirdecayslope*(*x*, *meth*) specifies the method for slope estimation. Possible values for *meth* are:
  - *meth* = ***Diff*** computes the slope as a ratio between the magnitude difference at the beginning and the ending of the decay period, and the corresponding time difference. (Default choice)
  - *meth* = ***Gauss*** computes the average of the slope, weighted by a gaussian curve that emphasizes values at the middle of the decay period (similar to Peeters, 2004).

## *mirduration*

### DESCRIPTION

Duration (in s.) of each successive event. It is computed by detecting the attack and decay phases of each event and by taking the portion of the curve between the onset and offset times that is over 40% of the maximum of the amplitude between the attack and decay times (Peeters, 2004).

### FLOWCHART INTERCONNECTIONS

*mirduration* accepts as input data type either:

- event detection curves (resulting from *mirevents*),
- and all the input data accepted by *mirevents*. In this case *mirevents* is called with the '**Filter**' method.

*mirevents* is normalized using the '**Normal**' option set to '**AcrossSegments**'.

Some options in *mirevents* can be controlled:

- *mirduration*(..., '**Single**') toggles on the '*Single*' option in *mirevents*.
- *mirduration*(..., '**LogCurve**') toggles on the '*Log*' option in *mirevents*.
- *mirduration*(..., '**MinLog**', *ml*) controls the '*MinLog*' option in *mirevents*.
- *mirduration*(..., '**PreSilence**', 0) toggles off the '*PreSilence*' option in *mirevents*.
- *mirduration*(..., '**PostSilence**') toggles on the '*PostSilence*' option in *mirevents*

The peak picking from the event detection is performed in any case. Its '**Contrast**' parameter can be specified. Its default value is the same as in *mirevents*.

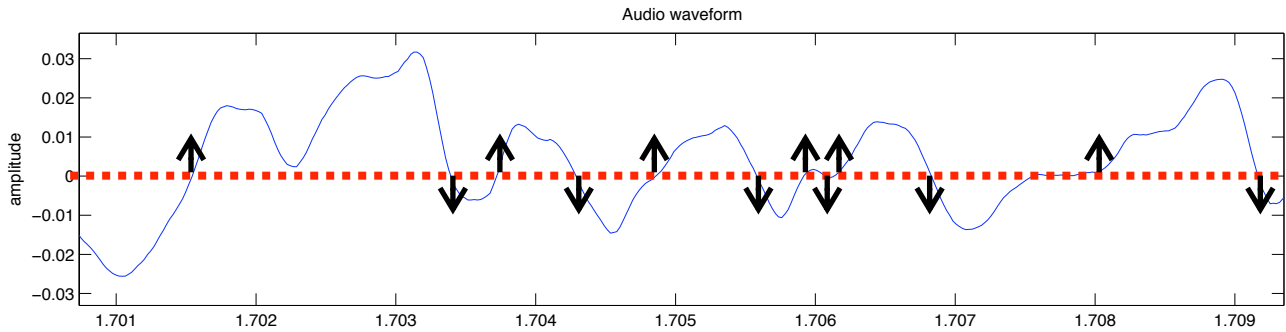
*mirduration* can return several outputs:

1. the duration itself and
2. the event detection curve returned by *mirevents*, including the detected events and the attack and decay phases.

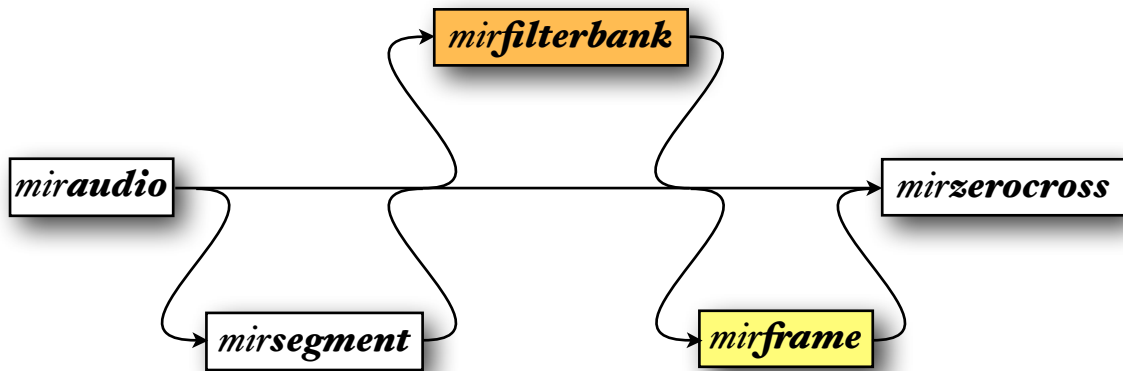
## *mirzerocross*

### WAVEFORM SIGN-CHANGE RATE

A simple indicator of noisiness consists in counting the number of times the signal crosses the X-axis (or, in other words, changes sign).



### FLOWCHART INTERCONNECTIONS



*mirzerocross* actually accepts any input data type (cf. section 4.2).

### FRAME DECOMPOSITION

*mirzerocross*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

### OPTIONS

- *mirzerocross*(..., **'Per'**, *p*) precises the temporal reference for the rate computation. Possible values:
  - *p* = *'Second'*: number of sign-changes per second (Default).

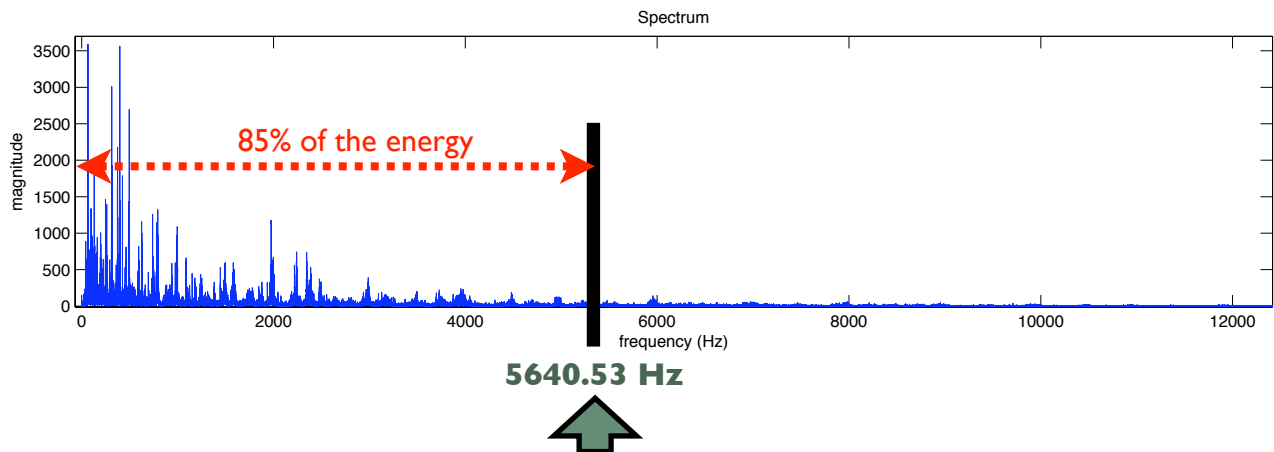


- $p = \text{'Sample'}$ : number of sign-changes divided by the total number of samples. The  $\text{'Second'}$  option returns a result equal to the one returned by the  $\text{'Sample'}$  option multiplied by the sampling rate.
- $\text{mirzerocross}(\dots, \text{'Dir'}, d)$  precises the definition of sign change. Possible values:
  - $d = \text{'One'}$ : number of sign-changes from negative to positive only (or, equivalently, from positive to negative only). (Default)
  - $d = \text{'Both'}$ : number of sign-changes in both ways. The  $\text{'Both'}$  option returns a result equal to twice the one returned by the  $\text{'One'}$  option.

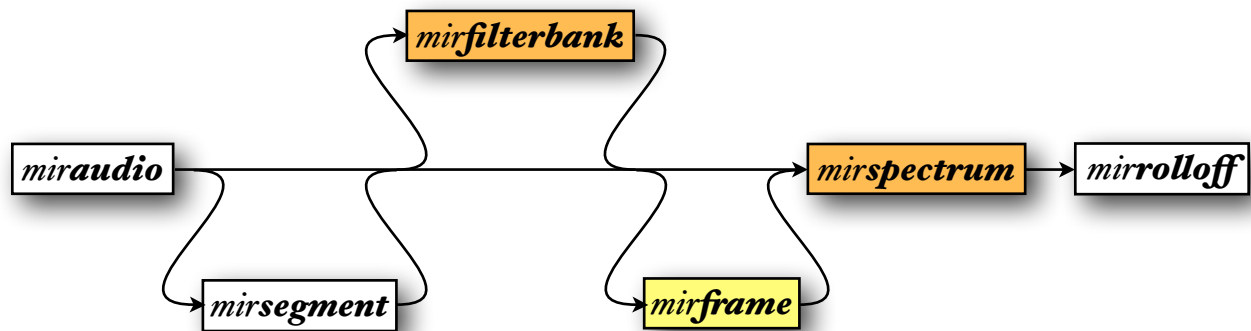
## *mirrolloff*

### HIGH-FREQUENCY ENERGY (I)

One way to estimate the amount of high frequency in the signal consists in finding the frequency such that a certain fraction of the total energy is contained below that frequency. This ratio is fixed by default to .85 (following Tzanetakis and Cook, 2002), other have proposed .95 (Pohle, Pampalk and Widmer, 2005).



### FLOWCHART INTERCONNECTIONS



*mirrolloff* accepts either:

- *mirspectrum* objects, or
- *miraudio* objects (same as for *mirspectrum*),
- **file name** or the **'Folder'** keyword.

### FRAME DECOMPOSITION

*mirrolloff*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

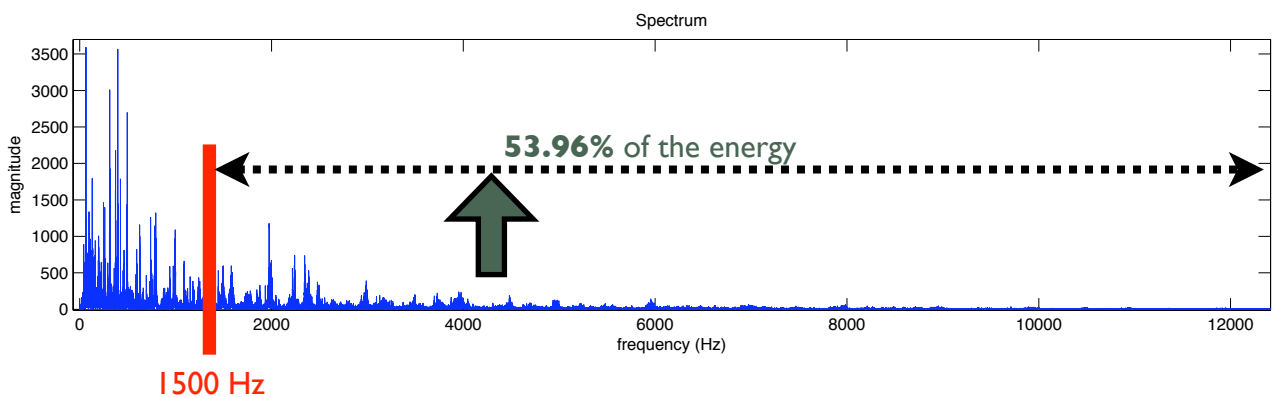
## OPTION

- *mirrolloff*(..., '**Threshold**', *p*) specifies the energy threshold, as a percentage. Default value: .85
- *mirrolloff*(..., '**MinRMS**', *m*) specifies the threshold *m*, as a value from 0 to 1, for the detection of quasi-silent frames, for which no value is given. For a given frame, if the RMS of the input spectrogram is below *m* times the highest RMS value over the frames, *NaN* is returned. default value: *m* = .005

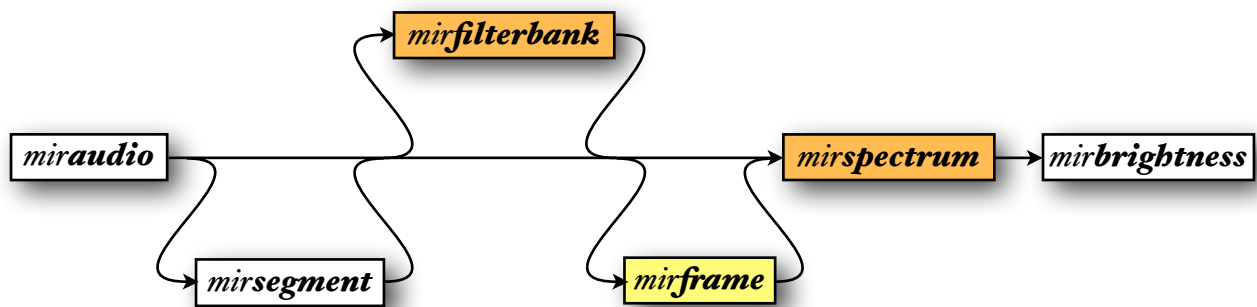
## *mirbrightness*

### HIGH-FREQUENCY ENERGY (II)

A dual method consists in fixing this time the cut-off frequency, and measuring the amount of energy above that frequency (Juslin, 2000). The result is expressed as a number between 0 and 1.



### FLOWCHART INTERCONNECTIONS



*mirbrightness* accepts either:

- ***mirspectrum*** objects, or
- ***miraudio*** objects (same as for *mirspectrum*).
- **file name** or the **'Folder'** keyword.

### FRAME DECOMPOSITION

*mirbrightness*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. **'Frame'** section.

## OPTIONS

- *mirbrightness*(..., '**CutOff**', *f*) specifies the frequency cut-off, in Hz. Default value: 1500 Hz. The value 1000 Hz has been proposed in (Laukka, Juslin and Bresin, 2005), and the value of 3000 Hz has been proposed in (Juslin, 2000).
- *mirbrightness*(..., '**MinRMS**', *m*) specifies the threshold *m*, as a value from 0 to 1, for the detection of quasi-silent frames, for which no value is given. For a given frame, if the RMS of the input spectrogram is below *m* times the highest RMS value over the frames, *NaN* is returned. default value: *m* = .005

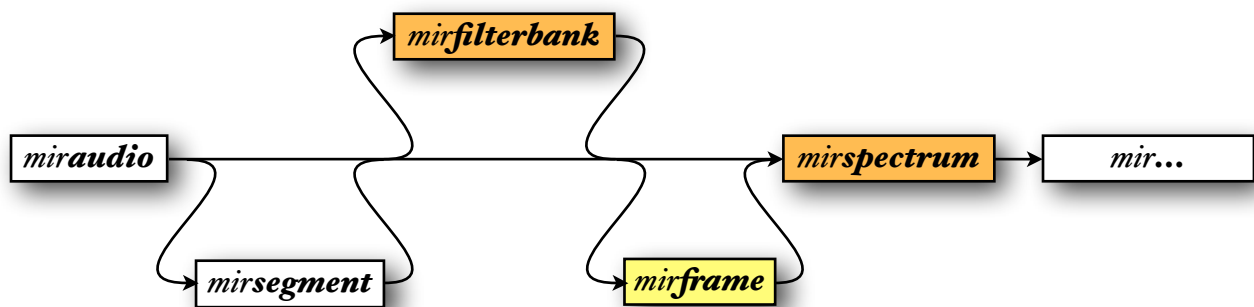
## *mircentroid, mirspread, mirskewness, mirkurtosis, mirflatness, mirentropy*

### STATISTICAL DESCRIPTION OF SPECTRAL DISTRIBUTION

The spectral distribution can be described by statistical moments: centroid, spread, skewness, kurtosis, flatness as well as by entropy.

See the *Statistics* section (§4.2) for more details about each operator.

### FLOWCHART INTERCONNECTIONS



All these operators accepts either:

- ***mirspectrum*** objects
- ***miraudio*** objects (same as for *mirspectrum*)
- **file name** or the ***Folder*** keyword.

### FRAME DECOMPOSITION

***Frame*** performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *Frame* section.

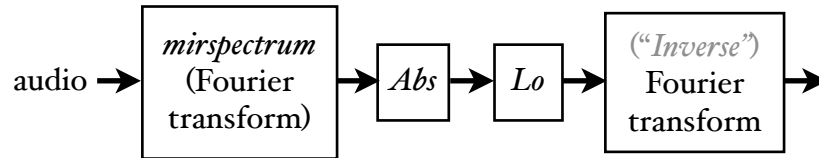
### OPTIONS

See the *Statistics* section (§4.2) for more details about each operator.

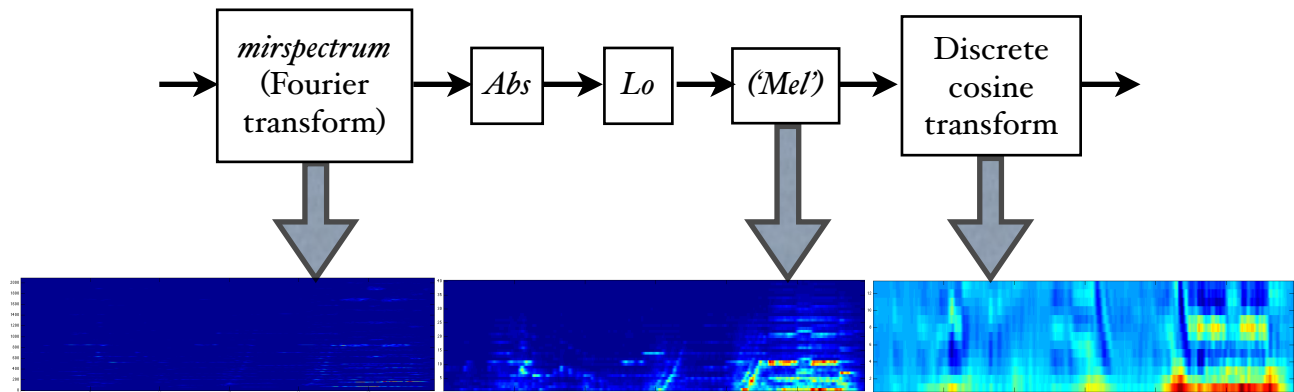
## *mirmfcc*

### MEL-FREQUENCY CEPSTRAL COEFFICIENTS

MFCC offers a description of the spectral shape of the sound. We recall that the computation of the cepstrum followed the following scheme:



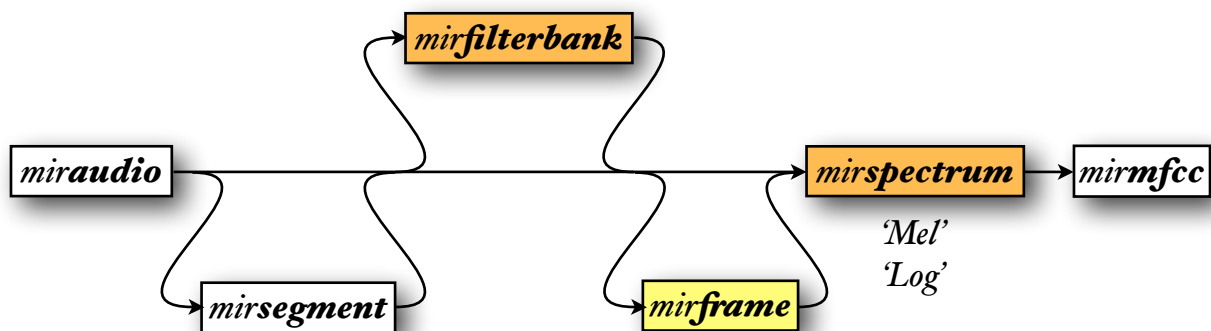
The computation of mel-frequency cepstral coefficients is highly similar:



Here the frequency bands are positioned logarithmically (on the Mel scale) which approximates the human auditory system's response more closely than the linearly-spaced frequency bands. And the Fourier Transform is replaced by a Discrete Cosine Transform. A discrete cosine transform (DCT) is a Fourier-related transform similar to the discrete Fourier transform (DFT), but using only real numbers. It has a strong "energy compaction" property: most of the signal information tends to be concentrated in a few low-frequency components of the DCT. That is why by default only the first 13 components are returned.

By convention, the coefficient of rank *zero* simply indicates the average energy of the signal.

### FLOWCHART INTERCONNECTIONS



*mirmfcc* accepts either:

- ***mirspectrum*** objects, or
- ***miraudio*** objects (same as for *mirspectrum*),
- **file name** or the **'Folder'** keyword.

*mirmfcc* can return several outputs:

- the mfcc coefficients themselves and
- the spectral representation (output of *mirspectrum*), in mel-band and log-scale.

## FRAME DECOMPOSITION

*mirmfcc*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

## OPTIONS

- *mirmfcc*(..., **'Bands'**, *b*) indicates the number of bands used in the mel-band spectrum decomposition. By default, *b* = 40.
- *mirmfcc*(..., **'Rank'**, *N*) computes the coefficients of rank(s) *N*. The default value is *N* = 1:13. Beware that the coefficient related to the average energy is by convention here of rank 0. This zero value can be included to the array *N* as well.
- If the output is frame-decomposed, showing the temporal evolution of the MFCC along the successive frames, the temporal differentiation can be computed:
  - *mirmfcc*(..., **'Delta'**, *d*) performs temporal differentiations of order *d* of the coefficients, also called delta-MFCC (for *d* = 1) or delta-delta-MFCC (for *d* = 2). By default, *d* = 1.
  - *mirmfcc*(..., **'Radius'**, *r*) specifies, for each frame, the number of successive and previous neighbouring frames taken into consideration for the least-square approximation used for the derivation. For a given radius *r*, the Delta operation for each frame *i* is computed by summing the MFCC coefficients at frame *i*+*j* (with *j* from -*r* to +*r*), each coefficient being multiplied by its weight *j*. Usually the radius is equal to 1 or 2. Default value: *r* = 2.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **'Rank'**: the series of rank(s) taken into consideration (same as *'Pos'*),

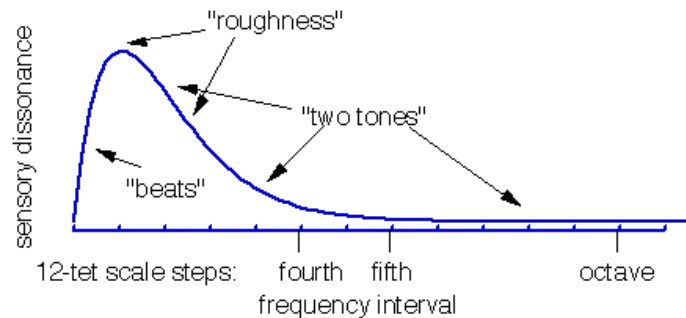


- ***Delta***: the number of times the delta operation has been performed.

## *mirroughness*

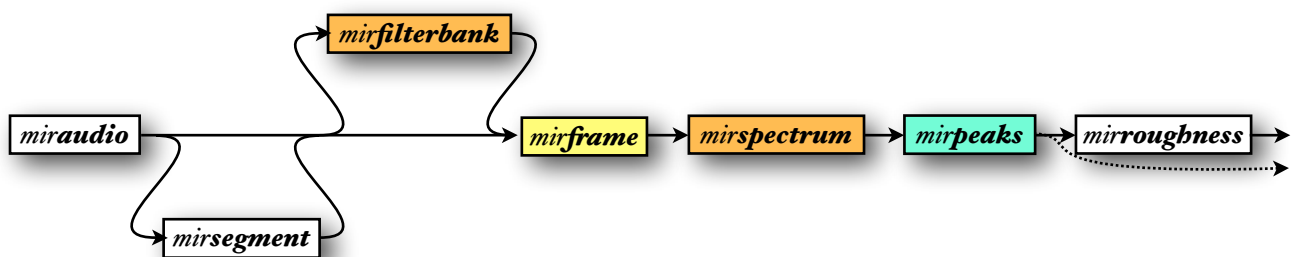
### SENSORY DISSONANCE

Plomp and Levelt (1965) have proposed an estimation of the sensory dissonance, or roughness, related to the beating phenomenon whenever pair of sinusoids are closed in frequency. The authors propose as a result an estimation of roughness depending on the frequency ratio of each pair of sinusoids represented as follows:



An estimation of the total roughness is available in *mirroughness* by computing the peaks of the spectrum, and taking the average of all the dissonance between all possible pairs of peaks (Sethares, 1998).

### FLOWCHART INTERCONNECTIONS



The ‘**Contrast**’ parameter associated to *mirpeaks* can be specified, and is set by default to .01. *mirroughness* accepts either:

- *mirspectrum* objects, where peaks have already been picked or not,
- *miraudio* objects: same as for *mirspectrum*, except that a **frame decomposition** is automatically performed. This forced frame decomposition is due to the fact that roughness can only be associated to a spectral representation association to a short-term sound excerpt: there is no sensory dissonance provoked by a pair of sinusoid significantly distant in time.
- **file name** or the ‘**Folder**’ keyword.

*mirroughness* can return several outputs:

1. the roughness value itself and
2. the spectral representation (output of *mirspectrum*) showing the picked peaks (returned by *mirpeaks*).

## FRAME DECOMPOSITION

*mirroughness*(..., **'Frame'**, ...) specifies the frame configuration, with by default a frame length of 50 ms and half overlapping. For the syntax, cf. the previous *mirframe* vs. *'Frame'* section.

## OPTIONS

- *mirroughness*(..., *m*) specifies the method used:
  - *m* = **'Sethares'** (default): based on the summation of roughness between all pairs of sines (obtained through spectral peak-picking) (Sethares, 1998). For each pair of peaks of frequencies  $f_1$  and  $f_2$ , the corresponding elementary roughness is obtained by multiplying the two peak amplitudes  $m_1$  and  $m_2$  altogether, and by weighting the results with the corresponding factor  $d(f_1, f_2)$  given on the dissonance curve.

$$R = \sum_{\text{All pairs of peaks}} m_1 \times m_2 \times d(f_1, f_2)$$

- *mirroughness*(..., **'Min'**): Variant of the Sethares model where the summation is weighted by the minimum amplitude of each pair of peak, instead of the product of their amplitudes (Weisser and Lartillot, 2013).

$$R = \sum_{\text{All pairs of peaks}} \min(m_1, m_2) \times d(f_1, f_2)$$

- *m* = **'Vassilakis'**: variant of **'Sethares'** model with a more complex weighting (Vassilakis, 2001, Eq. 6.23).

$$R = \sum_{\text{All pairs of peaks}} m_1^{0.1} m_2^{0.1} \frac{\min(m_1, m_2)^{3.11}}{m_1 + m_2} d(f_1, f_2)$$

- *mirroughness*(..., **'Normal'**) normalises with respect to overall spectrum magnitude (Lartillot and Weisser, 2021). More precisely:
  - For the default Sethares model, the roughness computed on each frame is divided by the sum over frequency bins of the spectrum amplitudes  $m$ , squared (since the pair of magnitudes are also multiplied together on the numerator as well).

$$R = \frac{\sum_{\text{All pairs of peaks}} m_1 \times m_2 \times d(f_1, f_2)}{\sum m^2}$$

- For the Sethares model with ‘*Min*’ variant, the roughness computed on each frame is divided by the sum over frequency bins of the spectrum amplitudes  $m$ .

$$R = \frac{\sum_{\text{All pairs of peaks}} \min(m_1, m_2) \times d(f_1, f_2)}{\sum m}$$

- For the Vassilakis model, the roughness computed on each frame is divided by the sum over frequency bins of the spectrum amplitudes  $m$ , raised to the power of 0.2.

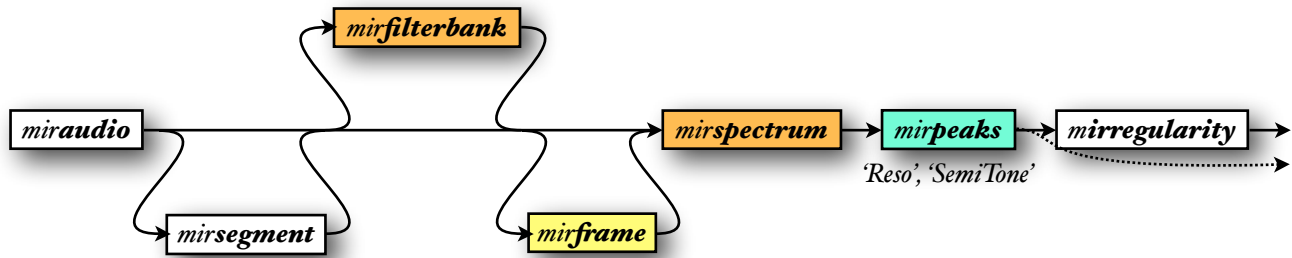
Whatever the chosen roughness model, in order to avoid getting high normalised roughness values for near-silent frames, these values should not be taken into account. For that reason, for any frame  $i$  for which the summation of spectrum amplitudes  $S_i$  is lower than a threshold value, the roughness value is replaced with a NaN value. The threshold value is defined as a ratio  $th$  (between 0 and 1) multiplied with the maximum of  $S_i$  across all frames  $i$ . The default value for  $th$  is 0.1, and can be specified using the command `mirroughness(..., ‘NormalThreshold’,  $th$ )`

## mirregularity

### SPECTRAL PEAKS VARIABILITY

The irregularity of a spectrum is the degree of variation of the successive peaks of the spectrum.

### FLOWCHART INTERCONNECTIONS



The '**Contrast**' parameter associated to **mirpeaks** can be specified, and is set by default to .or  
*mirregularity* accepts either:

- **mirspectrum** objects, where peaks have already been picked or not,
- **miraudio** objects (same as for *mirspectrum*),
- **file name** or the '**Folder**' keyword.

### FRAME DECOMPOSITION

*mirregularity*(..., '**Frame**', ...) performs first a frame decomposition, with by default a frame length of 50 ms and half overlapping. For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. '**Frame**' section.

### OPTIONS

- *mirregularity*(..., '**Jensen**') is based on (Jensen, 1999), where the irregularity is the sum of the square of the difference in amplitude between adjoining partials.(Default approach)

$$\left( \sum_{k=1}^N (a_k - a_{k+1})^2 \right) / \sum_{k=1}^N a_k^2$$

- *mirregularity*(..., '**Krimphoff**') is based on (Krimphoff et al., 1994), where the irregularity is the sum of the amplitude minus the mean of the preceding, same and next amplitude.

$$\sum_{k=2}^{N-1} \left| a_k - \frac{a_{k-1} + a_k + a_{k+1}}{3} \right|$$

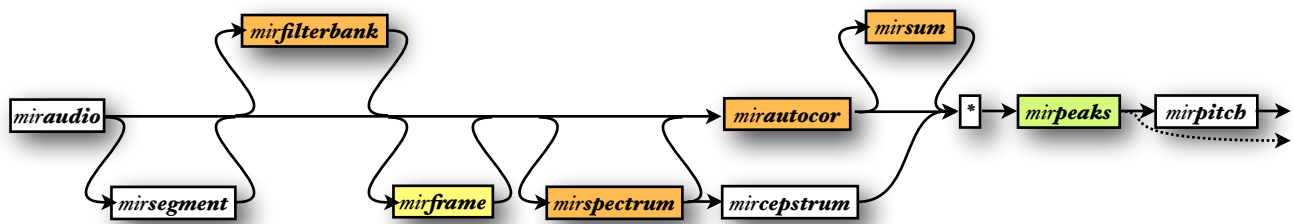
### 3.4. Pitch

## *mirpitch*

### PITCH ESTIMATION

Extract pitches, returned either as continuous pitch curves or as discretized note events.

### FLOWCHART INTERCONNECTIONS



The pitch content can be estimated in various ways:

- *mirpitch*(..., '**Autocor**') computes an autocorrelation function of the audio waveform, using *mirautocor*. This is the default strategy. Options related to *mirautocor* can be specified:
  - '**Enhanced**', toggled on by default here,
  - '**Compress**', set by default to .5,
  - filterbank configuration can be specified: either '**2Channels**' (default configuration), '**Gammatone**' or '**NoFilterbank**',
  - if a filterbank is used, '**Sum**' specifies whether the channels are recombined once the autocorrelation function is computed ('Sum', 1, which is the default), or if on the contrary, the channels are kept separate, and pitch content is extracted in each channel separately ('Sum', 0).
- *mirpitch*(..., '**Spectrum**') computes the FFT spectrum (*mirspectrum*). Options related to *mirspectrum* can be specified: '**Res**', '**dB**', '**Normal**' and '**Terhardt**'.
- *mirpitch*(..., '**AutocorSpectrum**') computes the autocorrelation (*mirautocor*) of the FFT spectrum (*mirspectrum*).
- *mirpitch*(..., '**Cepstrum**') computes the cepstrum (*mircepstrum*).
- These methods can be combined. In this case, the resulting representations (autocorrelation function or cepstral representations) are all expressed in the frequency domain and multiplied altogether.

Then a peak picking is applied to the autocorrelation function or to the cepstral representation. The parameters of the peak picking can be tuned.

- *mirpitch*(..., '**Total**', *m*) selects only the *m* best pitches.
- *mirpitch*(..., '**Mono**') only select the best pitch, corresponding hence to *mirpitch*(..., '**Total**', 1).
- *mirpitch*(..., '**Min**', *mi*) indicates the lowest pitch taken into consideration, in Hz. Default value: 75 Hz, following a convention in the *Praat* software (Boersma & Weenink, 2005).
- *mirpitch*(..., '**Max**', *ma*) indicates the highest pitch taken into consideration, expressed in Hz. Default value: 2400 Hz, because there seem to be some problems with higher frequency, due probably to the absence of pre-whitening in our implementation of Tolonen and Karjalainen autocorrelation approach (used by default).
- *mirpitch*(..., '**Threshold**', *t*) specifies the threshold factor for the peak picking. Default value: *c* = 0.4.
- *mirpitch*(..., '**Contrast**', *c*) specifies the contrast factor for the peak picking. Default value: *c* = 0.1.
- *mirpitch*(..., '**Order**', *o*) specifies the ordering for the peak picking. Default value: *o* = 'Amplitude'.
- *mirpitch*(..., '**Normalize**', *n*) specifies the normalisation for the peak picking. Default value: *n* = 'Global'.
- *mirpitch*(..., '**LocalFactor**', *f*) specifies the local factor for the peak picking. Default value when the keyword is mentioned: *f* = .99.

*mirpitch* accepts as input data type either:

- ***mirpitch*** objects,
- output of ***mirpeaks*** computation,
- ***mirautocor*** objects,
- ***mircepstrum*** objects,
- ***mirspectrum*** objects,
- ***miraudio*** objects, where the audio waveform can be:
  - segmented (using ***mirsegment***),



- when pitch is estimated by autocorrelating the audio waveform (*'Autocor'* strategy), the audio waveform is by default first decomposed into channels (cf. the ***Filterbank*** option below),
- decomposed into frames or not (using ***mirframe***);
- **file name** or the ***Folder*** keyword: same behavior than for *miraudio* objects,
- ***mirmidi*** objects.

*mirpitch* can return several outputs:

1. the pitch frequencies themselves, and
2. the ***mirautocor*** or ***mircepstrum*** data, where is highlighted the (set of) peak(s) corresponding to the estimated pitch (or set of pitches).

## FRAME DECOMPOSITION

*mirpitch*(..., ***Frame***, ...) performs first a frame decomposition, with by default a frame length of 46.4 ms and a hop factor of 10 ms (Tolonen & Karjalainen, 2000). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *Frame* section.

## POST-PROCESSING OPTIONS

- *mirpitch*(..., ***Cent***) convert the pitch axis from Hz to cent scale. One octave corresponds to 1200 cents, so that 100 cents correspond to a semitone in equal temperament.
- *mirpitch*(..., ***Segment***, *method*) segments the obtained monodic<sup>6</sup> pitch curve in cents as a succession of notes with stable frequencies.

Two methods are proposed:

- ***Lartillot*** (default method when ***Segment*** keyword is called). When this method is used for academic research, please cite the following publication:

Olivier Lartillot, “Computational analysis of maqam music: From audio transcription to musical analysis, everything is tightly intertwined”, *Acoustics 2012 Hong Kong Conference*.

Parameters used in this method:

---

<sup>6</sup> This ***Segment*** option requires a monodic pitch curve extraction using the ***Mono*** option, which is therefore toggled on, as well as the ***Cent*** and ***Frame*** options.

- *mirpitch*(..., '**SegMinLength**', *l*) specifies the minimum length of segments, in number of samples. Default length: *l* = 2 samples.
- *mirpitch*(..., '**SegPitchGap**', *g*) specifies the maximum tolerated pitch gap within a segment, in cents. The pitch distance is computed between the current sample is more distant and the average pitch of the segment before that point. Default gap: *g* = 45 cents.
- *mirpitch*(..., '**SegTimeGap**', *g*) specifies the maximum tolerated silence within a segment, in number of samples. Default gap: *g* = 20 samples.
- Segments with too low pitch periodicity (of autocorrelation value lower than 5% of the maximal value in that audio piece) are discarded.
- 'Nymoen' method from [pitchExtract.m](http://pitchExtract.m) by Kristian Nymoen. This method is particularly suitable for note segmentation of singing voice.
- *mirpitch*(..., '**Median**', *l*) performs a median filtering of the pitch curve. The length of the median filter is given by *l* (in s.). Its default value is .1 s. The median filtering can only be applied to mono-pitch curve. If several pitches were extracted in each frame, a mono-pitch curve is first computed by selecting the best peak of each successive frame.
- *mirpitch*(..., '**Stable**', *th*, *n*) remove pitch values when the difference (or more precisely absolute logarithmic quotient) with the *n* precedent frames exceeds the threshold *th*.
  - if *th* is not specified, the default value .1 is used.
  - if *n* is not specified, the default value 3 is used.
- *mirpitch*(..., '**Reso**', '**SemiTone**') removes peaks whose distance to one or several higher peaks is lower than a given threshold  $2^{(1/12)}$  (corresponding to a semitone).

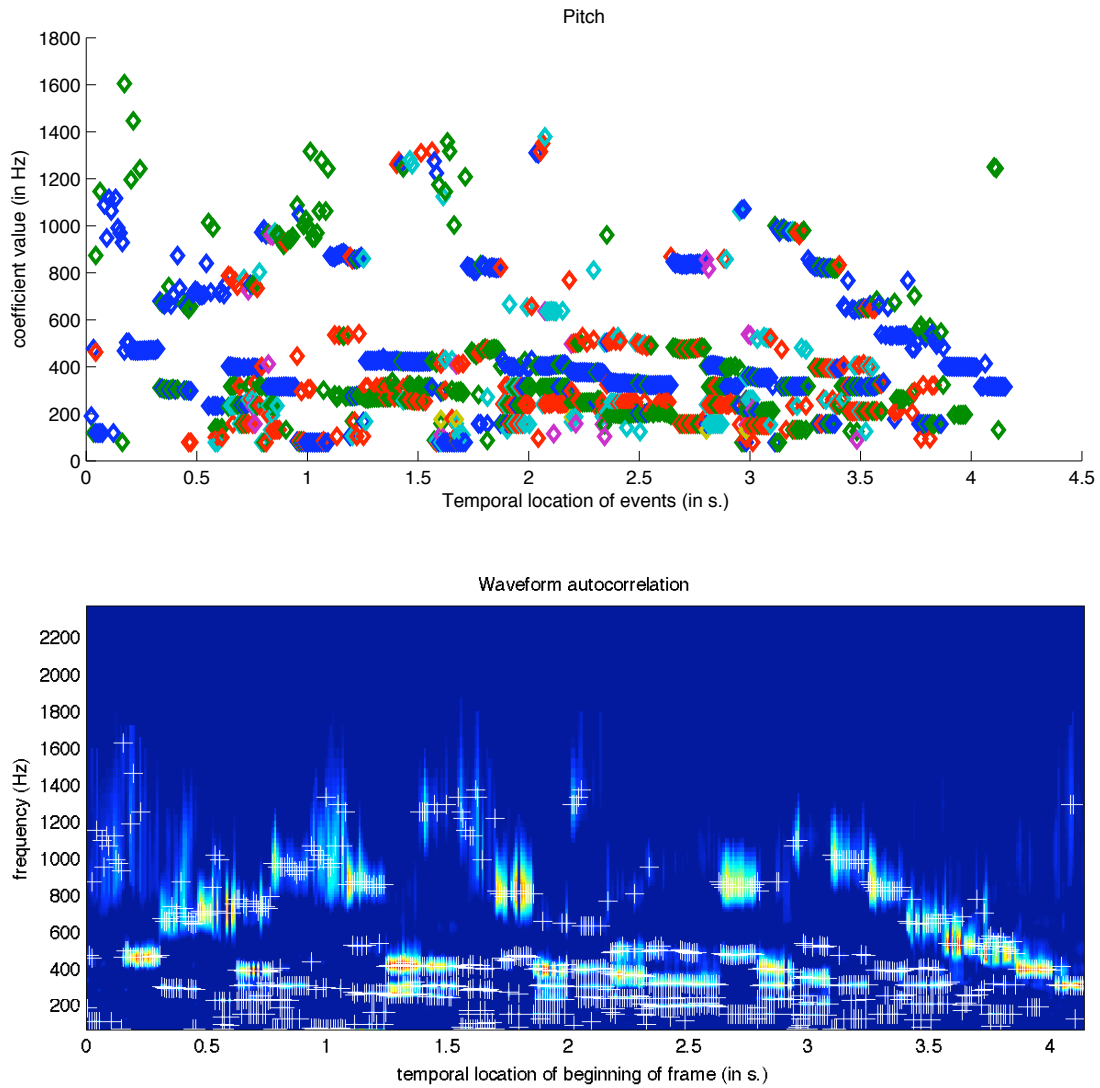
## P R E S E T   M O D E L

- *mirpitch*(..., '**Tolonen**') implements (part of) the model proposed in (Tolonen & Karjalainen, 2000). It is equivalent to

*mirpitch*(..., 'Enhanced', 2:10, 'Generalized', .67, '2Channels')

## E X A M P L E

$\{p_{ac}\} = \text{mirpitch}(\text{'ragtime.wav'}, \text{'Frame'})$



## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- ***Amplitude***: the amplitude associated with each pitch component.

## IMPORTATION OF PITCH DATA

*mirpitch*(*f*, *a*, *r*) creates a *mirpitch* object based on the frequencies specified in *f* and the related amplitudes specified in *a*, using a frame sampling rate of *r* Hz (set by default to 100 Hz).

Both *f* and *a* can be either:

- a matrix where each column represent one pitch track and lines corresponds to frames,
- an array of cells, where each cell, representing one individual frame, contains a vector.

## *mirmidi*

### AUTOMATED TRANSCRIPTION

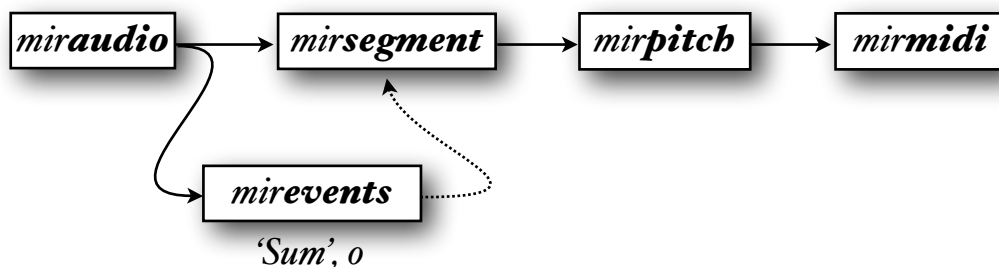
Segments the audio into events, extracts pitches related to each event and attempts a conversion of the result into a MIDI representation.

The audio segmentation is based on the event detection given by *mirevents* with the default parameter, but with the ‘*Sum*’ option toggled off in order to keep the channel decomposition of the input audio data.

The MIDI output is represented using the *MIDI Toolbox* note matrix representation. The displayed output is the piano-roll representation of the MIDI data, which requires *MIDI Toolbox*. Similarly, the result can be:

- sonified using *mirplay*, with the help of *MIDI Toolbox*;
- saved using *mirsave*:
  - as a MIDI file (by default), with the help of *MIDI Toolbox*,
  - as a *LilyPond* file (if the specified file has a ‘*.ly*’ extension).

### FLOWCHART INTERCONNECTIONS



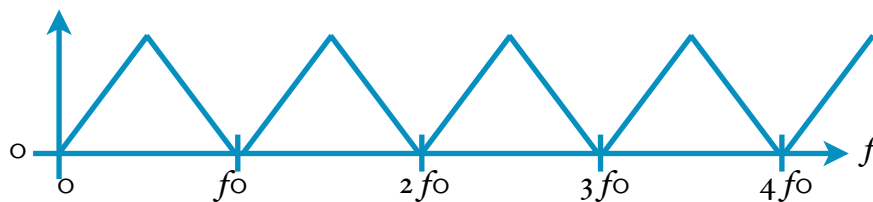
The ‘*Contrast*’ parameter associated to *mirpitch* can be specified, and is set by default to .3

## *mirinharmonicity*

### PARTIALS NON-MULTIPLE OF FUNDAMENTALS

*mirinharmonicity*(x) estimates the inharmonicity, i.e., the amount of partials that are not multiples of the fundamental frequency, as a value between 0 and 1. More precisely, the inharmonicity considered here takes into account the amount of energy outside the ideal harmonic series.

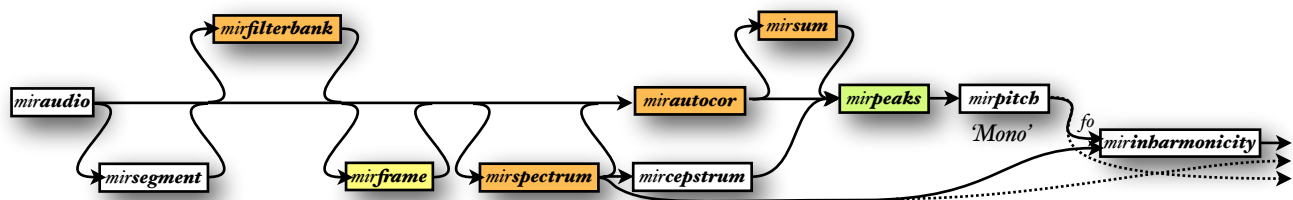
For that purpose, we use a simple function estimating the inharmonicity of each frequency given the fundamental frequency  $f_0$ :



#### **WARNING:**

This simple model presupposes that there is only one fundamental frequency.

### FLOWCHART INTERCONNECTIONS



*mirinharmonicity* accepts as main input either:

- ***mirspectrum*** objects,
- ***miraudio*** objects (same as for *mirspectrum*),
- **file name** or the **'Folder'** keyword.

*mirinharmonicity* can return several outputs:

1. the inharmonicity rate itself,
2. the ***mirspectrum*** data, and

3. the fundamental frequency '**fo**'.

## FRAME DECOMPOSITION

*mirinharmonicity*(..., '**Frame**', ...) performs first a frame decomposition, with by default a frame length of 10 ms and a hop factor of 12,5% (1.25 ms). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. '*Frame*' section.

## OPTION

*mirinharmonicity*(..., '**fo**', *f*) bases the computation of the inharmonicity on the fundamental frequency indicated by *f*. The frequency data can be either a number, or a *mirscalar* object (for instance, the output of a *mirpitch* computation).

By default, the fundamental frequency is computed using the command:

$$f = \text{mirpitch}(\dots, \text{'Mono'})$$

### 3.5. Tonality

#### *mirchromagram*

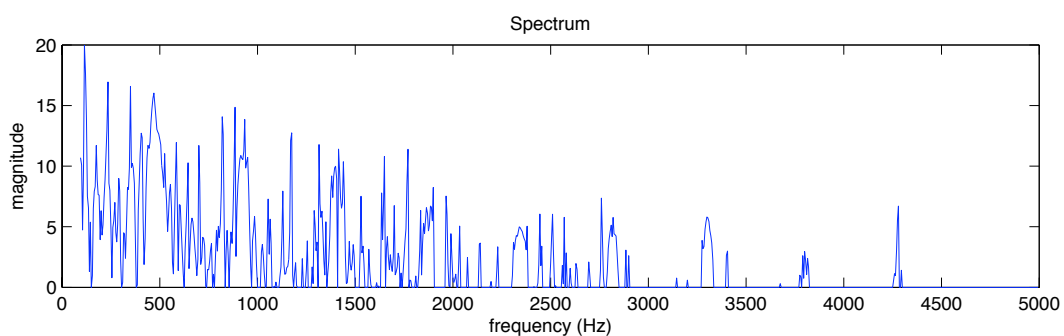
##### ENERGY DISTRIBUTION ALONG PITCHES

The chromagram, also called Harmonic Pitch Class Profile, shows the distribution of energy along the pitches or pitch classes.

- First the spectrum is computed in the logarithmic scale, with selection of, by default, the 20 highest dB, and restriction to a certain frequency range that covers an integer number of octaves, and normalization of the audio waveform before computation of the FFT.

$$s = \text{mirspectrum}(\dots, 'dB', 20, 'Min', fmin, 'Max', fmax, 'NormalInput', 'MinRes', r, 'OctaveRatio', .85)$$

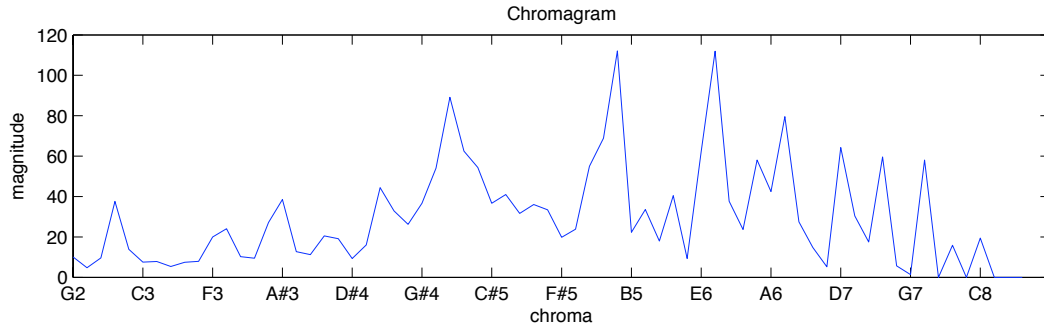
- The minimal frequency  $fmin$  is specified by the '*Min*' option (cf. below).
- The maximal frequency  $fmax$  is at least equal to the '*Max*' option, but the frequency range is extended if necessary in order to obtain an integer number of octaves.
- The minimal frequency resolution of the spectrum is chosen such that even lowest chromas will be segregated from the spectrum. It is based on the number of chromas per octave  $r$  ('*Res*' option, cf. below). Lowest frequencies are ignored if they do not meet this resolution constraint, and a warning message is sent by *mirspectrum* (cf. '*OctaveRatio*' keyword).



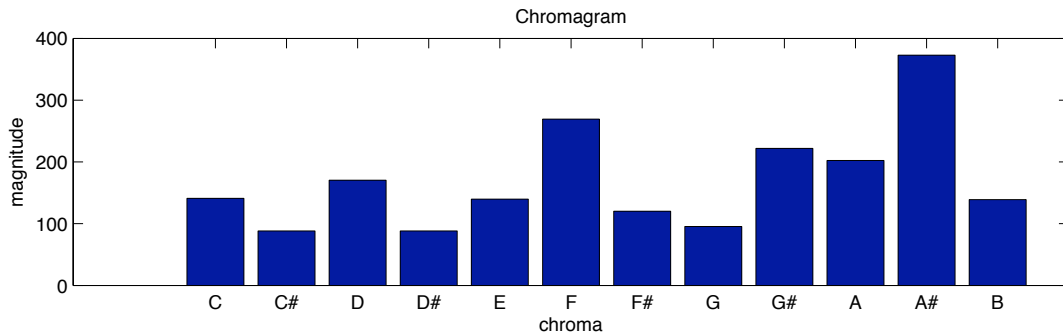
- The chromagram is a redistribution of the spectrum energy along the different pitches (i.e., “chromas”):

$$c = \text{mirchromagram}(s, 'Wrap', 'no')$$

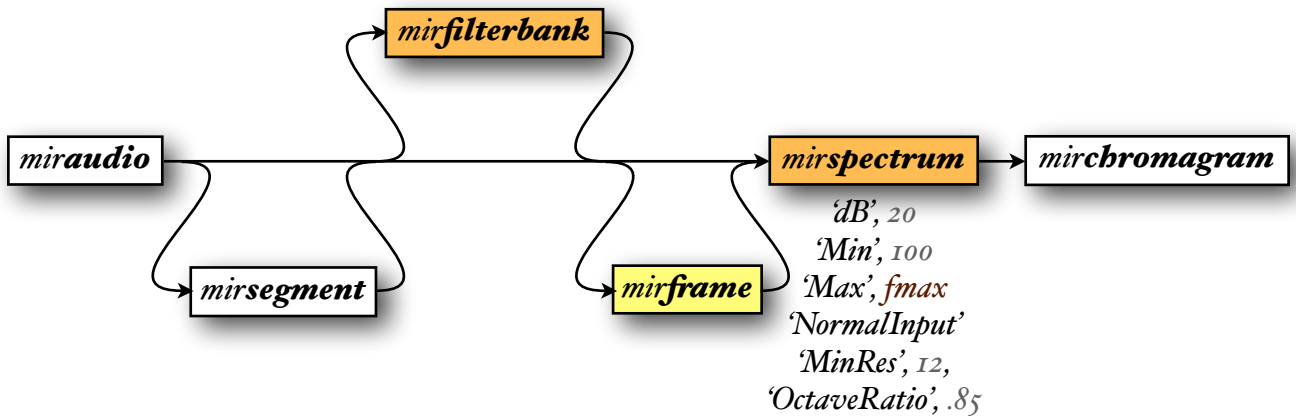
- If the '**Wrap**' option is selected, the representation is wrapped along the 12 pitch classes:



*c = mirchromagram(c, 'Wrap', 'yes')*



## FLOWCHART INTERCONNECTIONS



The '**Min**' and '**Max**' range used in *mirspectrum* can be tuned directly in *mirchromagram*, as well as the '**dB**' threshold (that can be written '**Threshold**' as well). These parameters are set by default to *Min* = 100 Hz, *Max* = 5000 Hz (Gómez, 2006) and *Threshold* = 20 dB. However, it seems that the spectrum should span as closely as possible an integer number of octaves, in order to avoid emphasizing particular pitches covered more often than others. Hence the higher limit of the frequency range of the spectrum computation is increased accordingly. Hence for the default parameters value (*Min* = 100 Hz, *Max* = 5000 Hz), the actual maximum frequency *fmax* is set to 6400 Hz. The '**MinRes**' value corresponds to the '*Res*' chromagram resolution parameter, as explained in the option section below.



*mirchromagram* accepts either:

- *mirspectrum* objects,
- *miraudio* objects, where the audio waveform can be segmented (using *mirsegment*), decomposed into channels (using *mirfilterbank*), and/or decomposed into frames (using *mirframe* or the **'Frame'** option, with by default a frame length of 200 ms and a hop factor of .05),
- **file name** or the **'Folder'** keyword.

## FRAME DECOMPOSITION

*mirchromagram*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 200 ms and a hop factor of 5% (10 ms). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. **'Frame'** section.

## OPTIONS

- $c = \text{mirchromagram}(\dots, \textbf{'Tuning'}, t)$  specifies the frequency (in Hz.) associated to chroma C. Default value,  $t = 261.6256$  Hz.
- $c = \text{mirchromagram}(\dots, \textbf{'Triangle'})$  weights the contribution of each frequency with respect to the distance with the actual frequency of the corresponding chroma.
- $c = \text{mirchromagram}(\dots, \textbf{'Weight'}, o)$  specifies the relative radius of the weighting window, with respect to the distance between frequencies of successive chromas.
  - $o = 1$ : each window begins at the centers of the previous one.
  - $o = .5$ : each window begins at the end of the previous one. (default value)
- $c = \text{mirchromagram}(\dots, \textbf{'Res'}, r)$  indicates the resolution of the chromagram in number of bins per octave. Default value,  $r = 12$ .

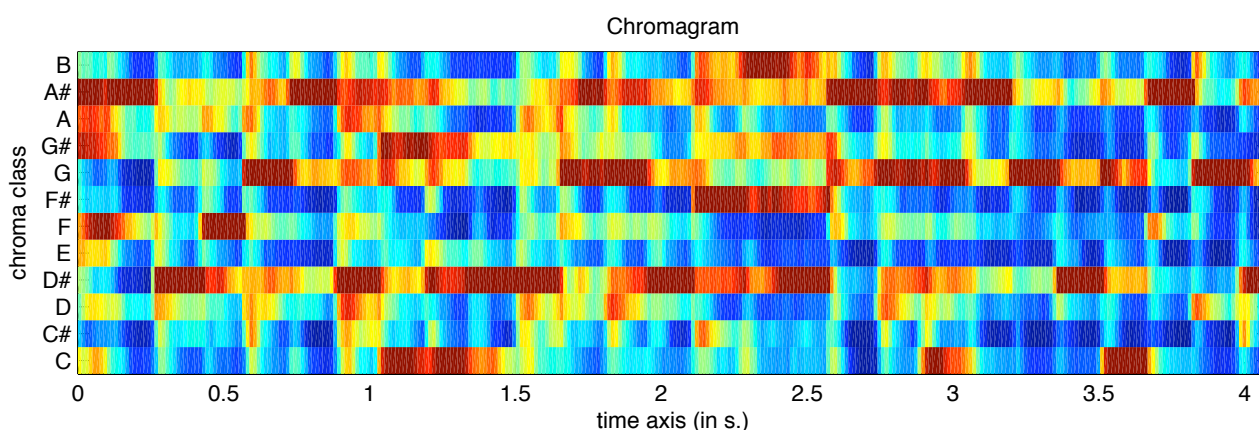
## POST-PROCESSING OPERATIONS

- $c = \text{mirchromagram}(\dots, \textbf{'Wrap'}, w)$  specifies whether the chromagram is wrapped or not.
  - $w = \text{'yes'}$ : groups all the pitches belonging to same pitch classes (default value)
  - $w = \text{'no'}$ : pitches are considered as absolute values.
- $c = \text{mirchromagram}(\dots, \textbf{'Center'})$  centers the result.

- $c = \text{mirchromagram}(\dots, \text{'Normal'}, n)$  normalizes the result, using the  $n$ -norm. The default value is  $n = \text{Inf}$ , corresponding to a normalization by the maximum value.  $n = 0$  toggles off the normalization. Alternative keyword: **'Norm'**.
- $c = \text{mirchromagram}(\dots, \text{'Pitch'}, p)$  specifies how to label chromas in the figures.
  - $p = \text{'yes'}$ : chromas are labeled using pitch names (default)
  - $p = \text{'no'}$ : chromas are labeled using MIDI pitch numbers.

## EXAMPLE

*mirchromagram('ragtime.wav', 'Frame')*



## ACCESSIBLE OUTPUT

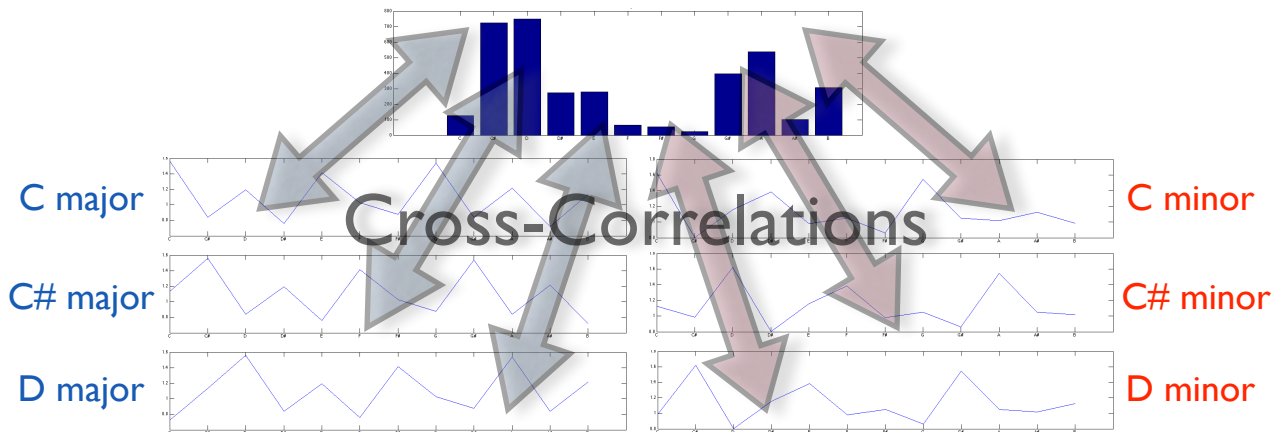
cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **'Magnitude'**: same as *'Data'*,
- **'Chroma'**: the chroma related to each magnitude (same as *'Pos'*),
- **'ChromaClass'**: the chroma class (A, A#, B, etc.) related to each chroma,
- **'ChromaFreq'**: the central frequency of each chroma, in the unwrapped representation,
- **'Register'**: the octave position,
- **'PitchLabel'**: whether the chroma are represented as simple numeric positions (o), or as couples (ChromaClass, Register) (i).
- **'Wrap'**: whether the chromagram is represented along all possible chromas (o), or along the 12 chroma-classes only (i).

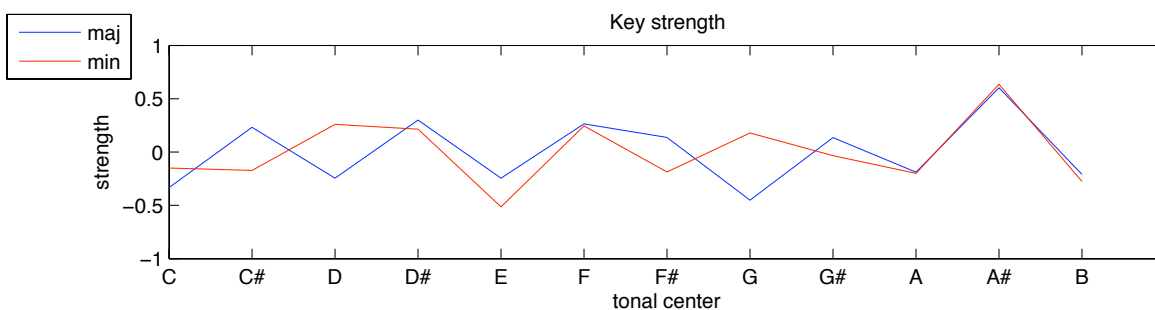
## mirkeystrength

### PROBABILITY OF KEY CANDIDATES

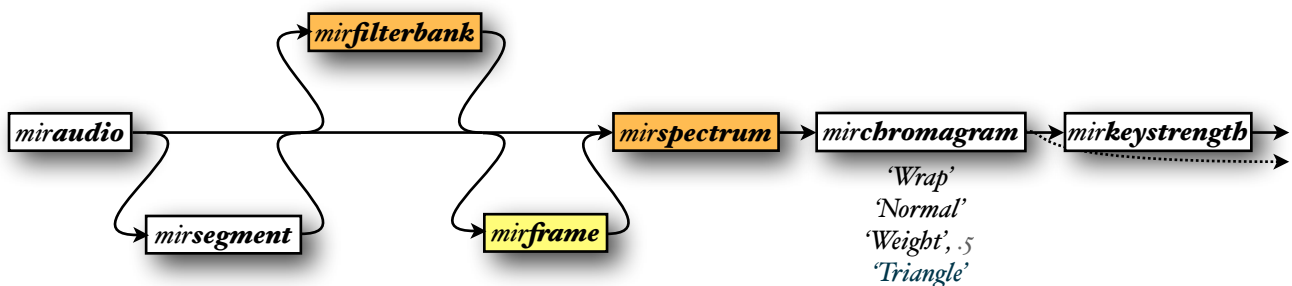
*mirkeystrength* computes the key strength, a score between -1 and +1 associated with each possible key candidate, by computing correlation coefficients between the chromagram returned by *mirchromagram*, wrapped and normalized (using the 'Normal' option), and similar profiles representing all the possible tonality candidates (Krumhansl, 1990; Gomez, 2006).



The resulting graph indicate the correlation coefficient for each different tonality candidate.



### FLOWCHART INTERCONNECTIONS



For the moment, only the '**Weight**' and '**Triangle**' options used in *mirchromagram* can be tuned directly in *mirkeystrength*.

*mirkeystrength* accepts either:

- *mirchromagram* objects,
- *mirspectrum* objects,
- *miraudio* objects (same as for *mirchromagram*),
- **file name** or the **'Folder'** keyword.

*mirkeystrength* can return several outputs:

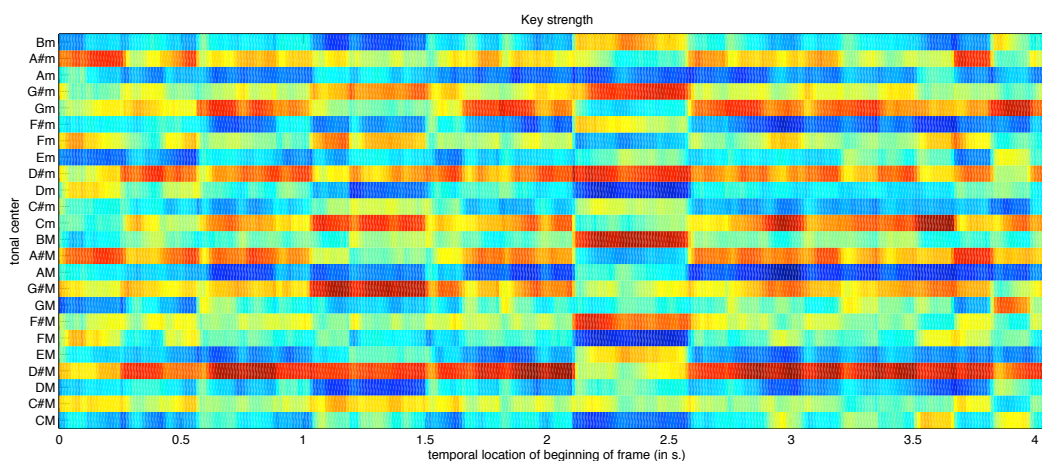
1. the key strength itself, and
2. the *mirchromagram* data.

## FRAME DECOMPOSITION

*mirkeystrength*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 200 ms and a hop factor of 5% (10 ms). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

## EXAMPLE

*mirkeystrength*('ragtime.wav', 'Frame')



## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

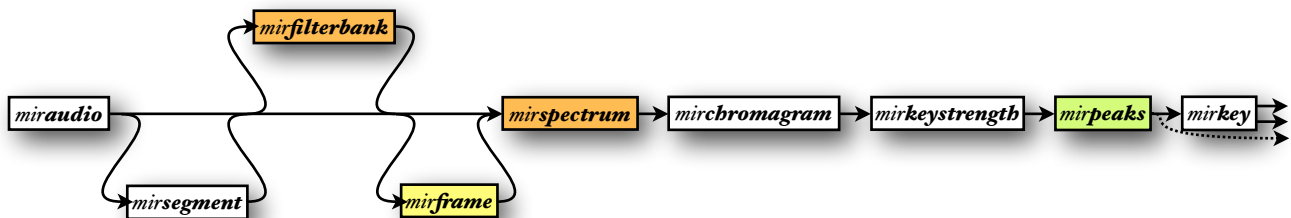
- **‘Strength’**: the key strength value for each key and each temporal position (same as *‘Data’*). The resulting matrix is distributed along two layers along the fourth dimension: a first layer for major keys, and a second layer for minor keys.
- **‘Tonic’**: the different key centres (same as *‘Pos’*).

## mirkey

### DESCRIPTION

Gives a broad estimation of tonal center positions and their respective clarity.

### FLOWCHART INTERCONNECTIONS



It consists simply of a peak picking in the *mirkeystrength* curve(s). Two options of *mirpeaks* are accessible from *mirkey*:

- **'Total'**, set to 1
- **'Contrast'**, set to .1

The **'Weight'** and **'Triangle'** options used in *mirchromagram* can be changed directly in *mirkeystrength*.

*mirkey* accepts either:

- *mirkeystrength* objects, where peaks have been already extracted or not,
- *mirchromagram* objects,
- *mirspectrum* objects,
- *miraudio* objects, where the audio waveform can be segmented (using *mirsegment*), decomposed into channels (using *mirfilterbank*), and/or decomposed into frames (using *mirframe* or the **'Frame'** option, with by default a frame length of 1 s and half overlapping),
- **file name** or the **'Folder'** keyword.

*mirkey* can return several outputs:

1. the best key(s), i.e., the peak abscissa(e);
2. the *key clarity*: i.e., the key strength associated to the best key(s), i.e., the peak ordinate(s);

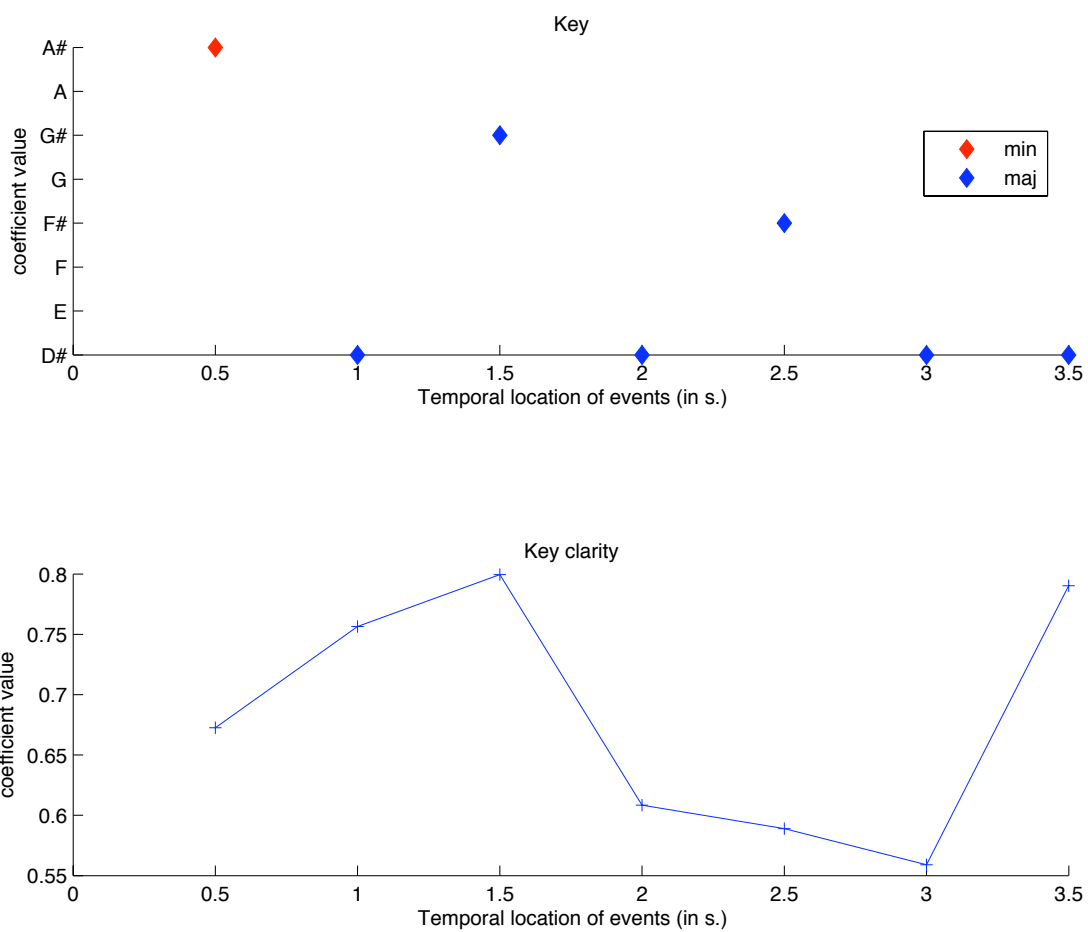
- the *mirkeystrength* data including the picked peaks (*mirpeaks*).

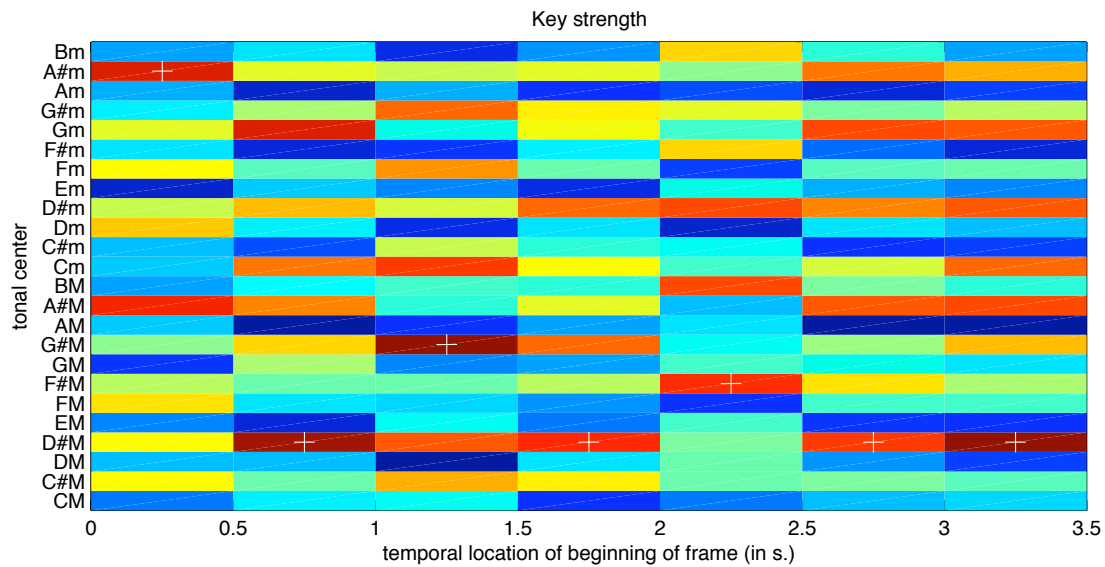
## FRAME DECOMPOSITION

*mirkey*(..., **'Frame'**, ...) performs first a frame decomposition, with by default a frame length of 1 s and a hop factor of 50% (0.5 s). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *'Frame'* section.

## EXAMPLE

$[k\ c\ s] = \text{mirkey}(\text{'ragtime.wav'}, \text{'Frame'})$



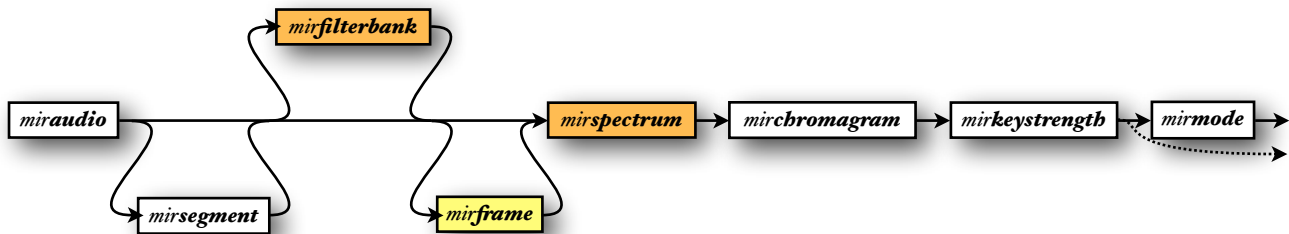


## *mirmode*

### DESCRIPTION

Estimates the modality, i.e. major vs. minor, returned as a numerical value between  $-1$  and  $+1$ : the closer it is to  $+1$ , the more major the given excerpt is predicted to be, the closer the value is to  $-1$ , the more minor the excerpt might be.

### FLOWCHART INTERCONNECTIONS



*mirmode* accepts either:

- ***mirkeystrength*** objects, where peaks have been already extracted or not,
- ***mirchromagram*** objects,
- ***mirspectrum*** objects,
- ***miraudio*** objects (same as for *mirkey*) or
- **file name** or the **'Folder'** keyword.



*mirmode* can return several outputs:

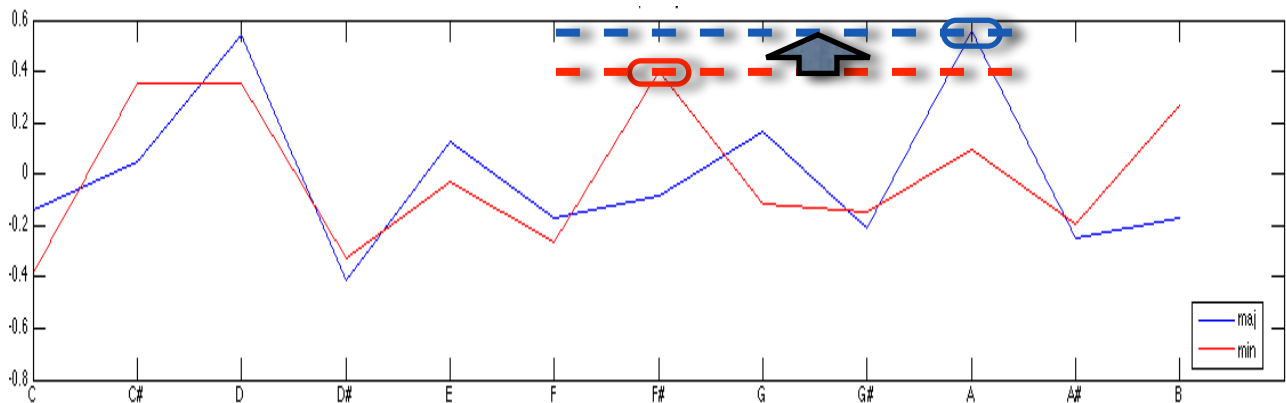
1. modality itself, and
2. the *mirkeystrength* result used for the computation of modality.

## FRAME DECOMPOSITION

*mirmode*(..., '**Frame**', ...) performs first a frame decomposition, with by default a frame length of 1 s and a hop factor of 50% (0.5 s). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. '*Frame*' section.

## STRATEGIES

- *mirmode*(..., '**Best**') computes the key strength difference between the best major key (highest key strength value) and the best minor key (lowest key strength value). (default choice)



- *mirmode*(..., '**Sum**') sums up the key strength differences between all the major keys and their relative minor keys.

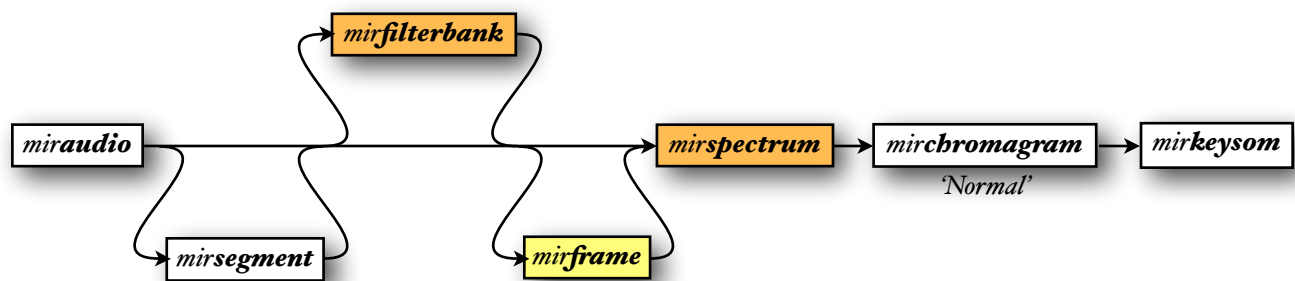
## mirkeysom

### DESCRIPTION

Projects the chromagram (normalized using the ‘*Normal*’ option) into a self-organizing map trained with the Krumhansl-Kessler profiles (modified for chromagrams) (Toivainen and Krumhansl, 2003; Krumhansl, 1990).

The result is displayed as a pseudo-color map, where colors correspond to Pearson correlation values. In case of frame decomposition, the projection maps are shown one after the other in an animated figure.

### FLOWCHART INTERCONNECTIONS



*mirkeysom* accepts either:

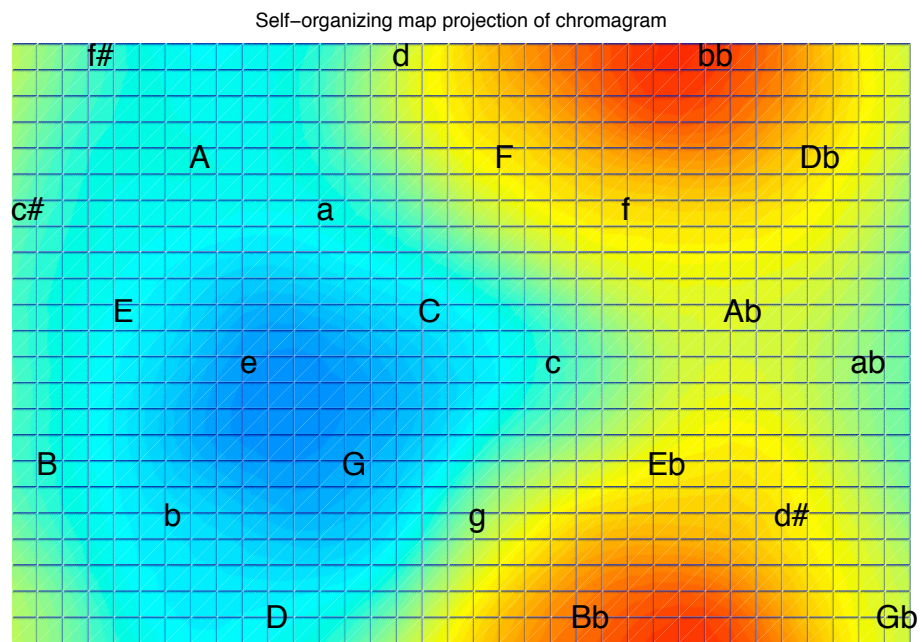
- *mirchromagram* objects,
- *mirspectrum* objects,
- *miraudio* objects , where the audio waveform can be segmented (using *mirsegment*), decomposed into channels (using *mirfilterbank*), and/or decomposed into frames (using *mirframe*),
- **file name** or the **‘Folder’** keyword.

### FRAME DECOMPOSITION

*mirkeysom*(..., **‘Frame’**, ...) performs first a frame decomposition, with by default a frame length of 1 s and a hop factor of 50% (0.5 s). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. ‘*Frame*’ section.

### EXAMPLE

*mirkeysom*(‘ragtime.wav’)



## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

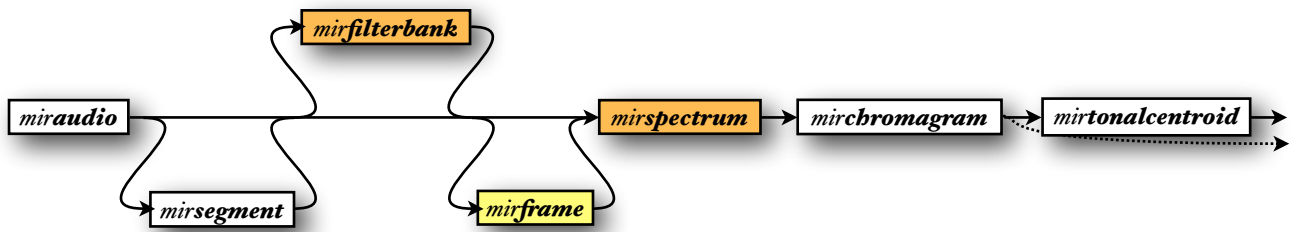
- ***Weight***: the projection value associated to each map location (same as '*Data*').

## *mirtonalcentroid*

### DESCRIPTION

Calculates the 6-dimensional tonal centroid vector from the chromagram. It corresponds to a projection of the chords along circles of fifths, of minor thirds, and of major thirds (Harte and Sandler, 2006).

### FLOWCHART INTERCONNECTIONS



*mirtonalcentroid* accepts either:

- ***mirchromagram*** objects,
- ***mirspectrum*** objects,
- ***miraudio*** objects, where the audio waveform can be segmented (using ***mirsegment***), decomposed into channels (using ***mirfilterbank***), and/or decomposed into frames (using ***mirframe***),
- **file name** or the ***Folder*** keyword.

*mirtonalcentroid* can return several outputs:

1. the tonal centroid itself, and
2. the ***mirchromagram*** data.

### FRAME DECOMPOSITION

*mirtonalcentroid*(..., ***Frame***, ...) performs first a frame decomposition, with by default a frame length of 743 ms and a hop factor of 10% (74.3 ms). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *Frame* section.

### ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- ***Dimensions***: the index of the 6 dimensions (same as *Pos*),

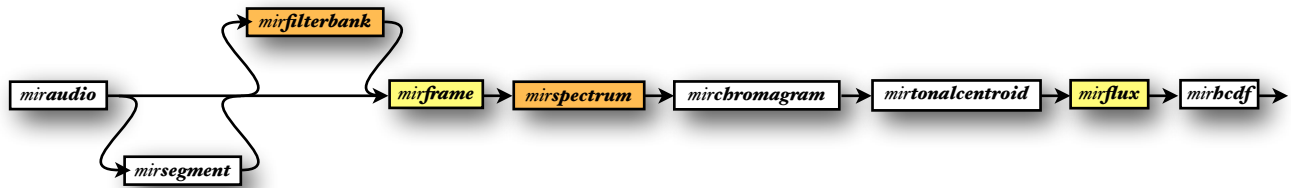
- ***Positions***: the position of each data within the 6-dimensional space (same as *Data*).

## *mirhcdf*

### DESCRIPTION

The Harmonic Change Detection Function (HCDF) is the flux of the tonal centroid (Harte and Sandler, 2006).

### FLOWCHART INTERCONNECTIONS



*mirhcdf* accepts either:

- ***mirtonalcentroid*** frame-decomposed objects,
- ***mirchromagram*** frame-decomposed objects,
- ***mirspectrum*** frame-decomposed objects,
- ***miraudio*** objects , where the audio waveform can be segmented (using ***mirsegment***), decomposed into channels (using ***mirfilterbank***). If not decomposed yet, it is decomposed into frames (using the ***Frame*** option, with by default a frame length of .743 s and a hop factor of .1),
- **file name** or the ***Folder*** keyword.

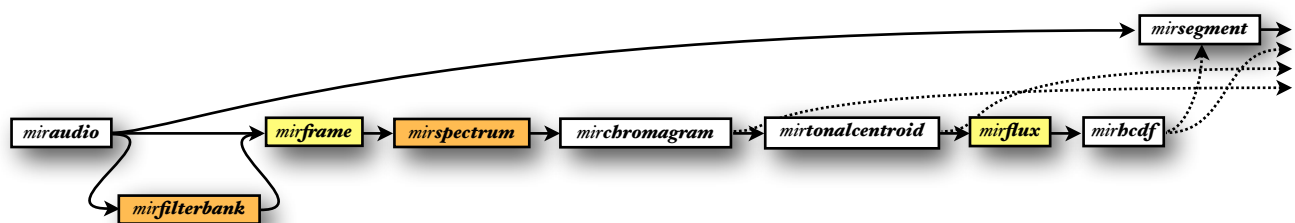
### FRAME DECOMPOSITION

*mirhcdf*(..., ***Frame***, ...) performs first a frame decomposition, with by default a frame length of 743 ms and a hop factor of 10% (74.3 ms). For the specification of other frame configuration using additional parameters, cf. the previous *mirframe* vs. *Frame* section.

## *mirsegment(..., 'HCDF')*

Peak detection applied to the HCDF returns the temporal position of tonal discontinuities that can be used for the actual segmentation of the audio sequence.

### FLOWCHART INTERCONNECTIONS



*mirsegment* accepts uniquely as main input a ***miraudio*** objects not frame-decomposed, not channel decomposed, and not already segmented. Alternatively, **file name** or the **'Folder'** keyword can be used as well.

*mirsegment(..., 'HCDF')* can return several outputs:

1. the segmented audio waveform itself,
2. the HCDF (***mirhcdf***) after peak picking (***mirpeaks***),
3. the tonal centroid (***mirtonalcentroid***), and
4. the chromagram (***mirchromagram***).

## 4. HIGH-LEVEL FEATURES

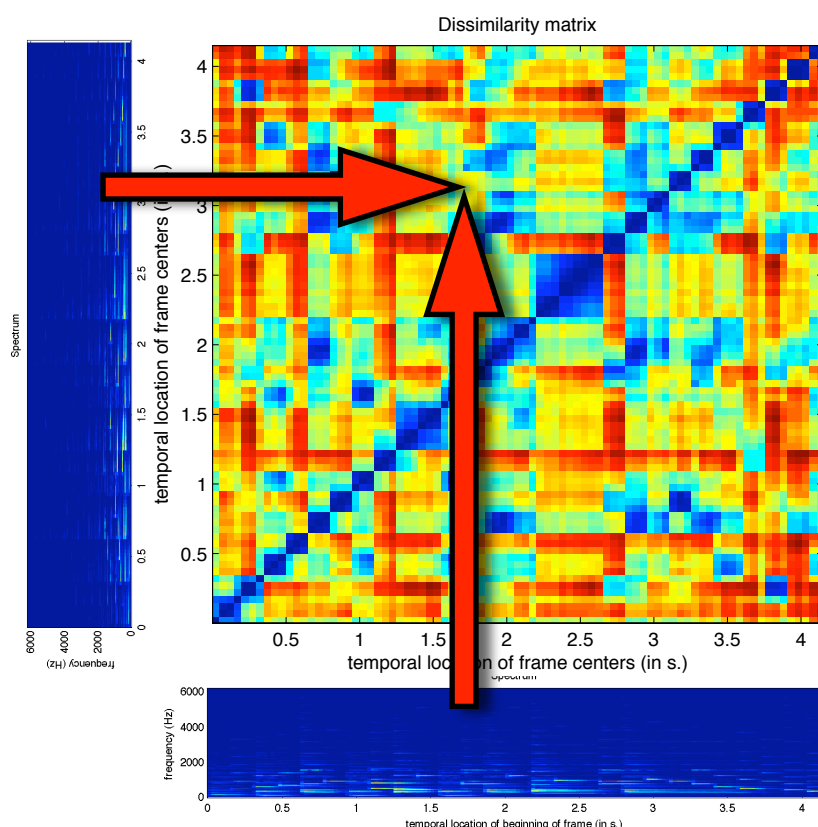
### 4.1. Structure and form

More elaborate tools have also been implemented that can carry out higher-level analyses and transformations. In particular, audio files can be automatically segmented into a series of homogeneous sections, through the estimation of temporal discontinuities along diverse alternative features such as timbre in particular (Foote and Cooper, 2003).

#### ***mirsimatrix***

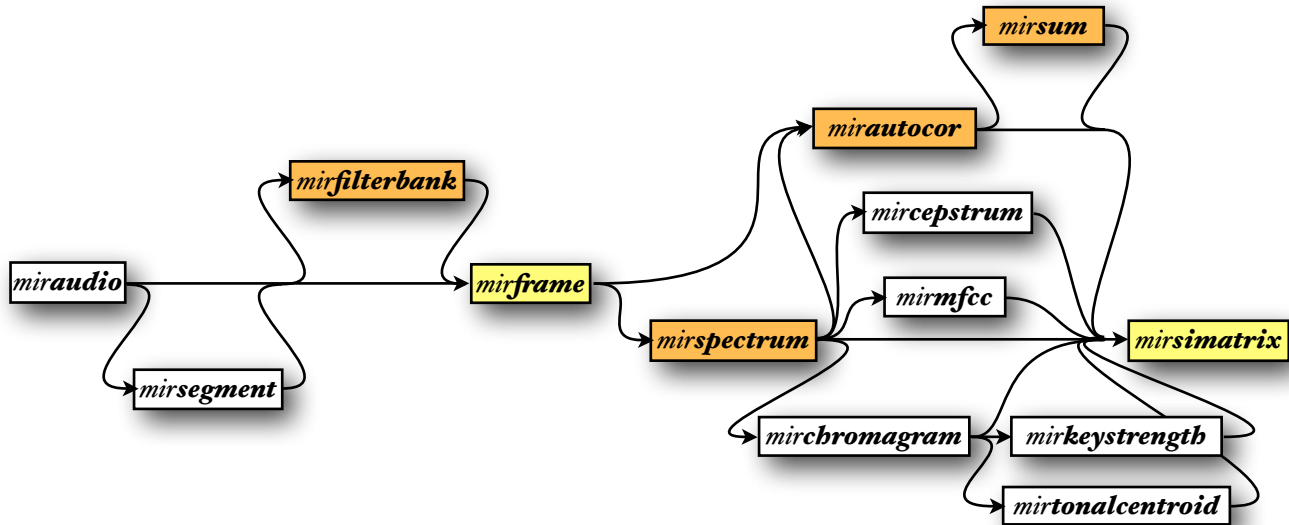
##### DESCRIPTION

A similarity matrix shows the similarity between all possible pairs of frames from the input data.





## FLOWCHART INTERCONNECTIONS



*mirsimatrix* usually accepts either:

- ***mirspectrum*** frame-decomposed objects.
- ***miraudio*** objects: in this case, the (dis)similarity matrix will be based on the spectrogram (***mirspectrum***). The audio waveform is decomposed into frames if it was not decomposed yet, and the default frame parameters – frame length of 50 ms and no overlapping – can be changed using the ***Frame*** option.

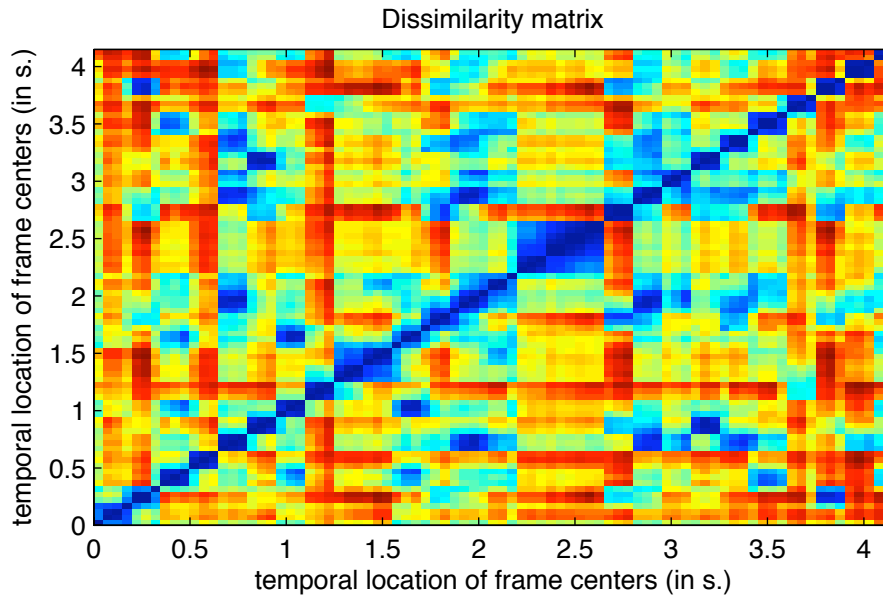
If the audio waveform is segmented (with ***mirsegment***), the similarities are not computed between frames but on the contrary between segments, using Kullback-Leibler distance (Foote and Cooper, 2003). (Removed since version 1.2, due to difficulties that have not been solved since.)

- **file name** or the ***Folder*** keyword: same behavior than for *miraudio* objects;
- ***mirautocor*** frame-decomposed objects;
- ***mircepstrum*** frame-decomposed objects;
- ***mirmfcc*** frame-decomposed objects;
- ***mirchromagram*** frame-decomposed objects;
- ***mirkeystrength*** frame-decomposed objects;
- a *Matlab* square matrix: in this case, a *mirsimatrix* object is created based on the matrix given as input, with the following options:
  - The frame rate used for the matrix can be specified as second input of *mirsimatrix*. By default, it is set to 20 Hz.

- The input is supposed to be a *similarity* matrix already. If the input is in fact a *dissimilarity* matrix, the keyword '*Dissimilarity*' should be added as third argument of *mirsimatrix*.

## OPTIONS

- *mirsimatrix*(..., '**Distance**', *f*) specifies the name of a distance function, from those proposed in the *Statistics Toolbox* (cf. *help pdist*). default value: *f* = '*cosine*'
- *mirsimatrix*(..., '**Width**', *w*) specifies the size of the diagonal bandwidth, in samples, outside which the dissimilarity will not be computed. If *w* is even, the actual width is *w*-1 samples. If *w* = *Inf* (default value), all the matrix will be computed. If the value is non-infinite<sup>7</sup>, a Hanning window is applied in order to smoothen the transition at the junction between the computed and non-computed regions of the matrix.
- *mirsimatrix*(..., '**Dissimilarity**') return the *dissimilarity matrix*, which is the intermediary result before the computation of the similarity matrix. It shows the distance between each possible frame.



- *mirsimatrix*(..., '**Similarity**', *f*) indicates the function *f* specifying the way the distance values in the dissimilarity matrix are transformed into the values of the *similarity matrix*. Two functions available:
  - *f* = '*oneminus*' (default parameter) corresponding to

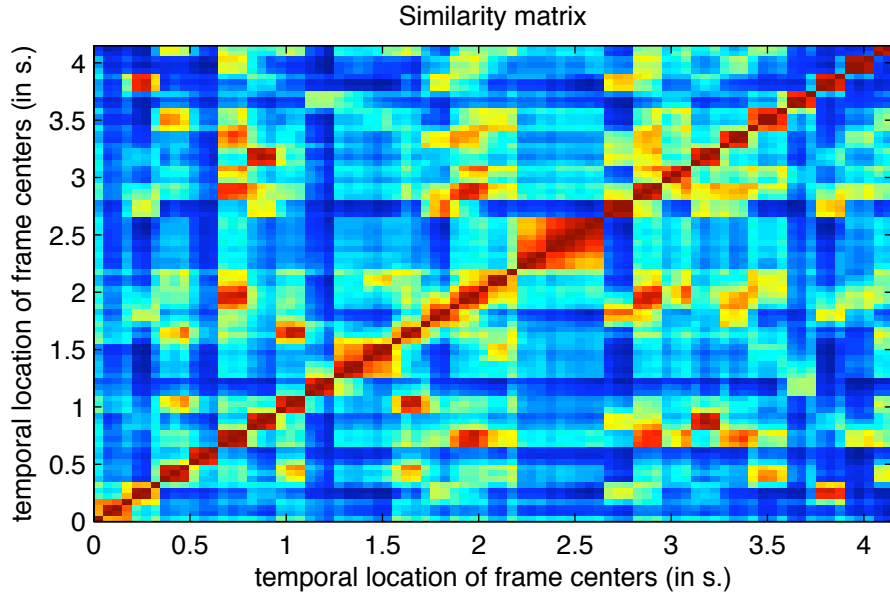
$$f(x) = 1 - x$$

<sup>7</sup> and if no '*Horizontal*' or '*TimeLag*' transformation is specified.

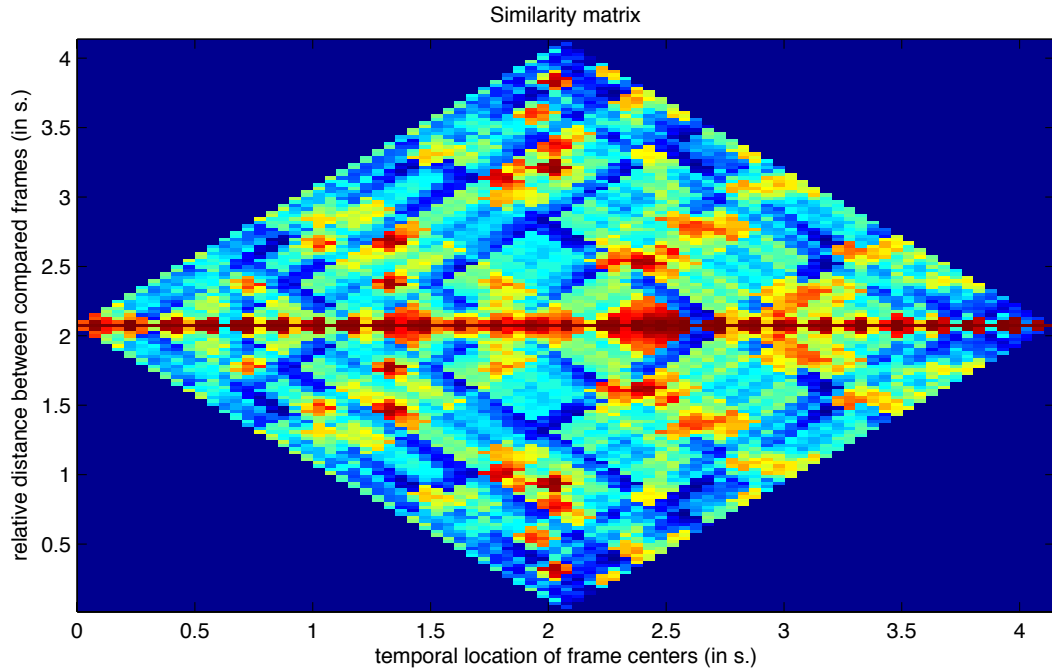
- $f = \text{'exponential'}$  corresponding to

$$f(x) = \exp(-x)$$

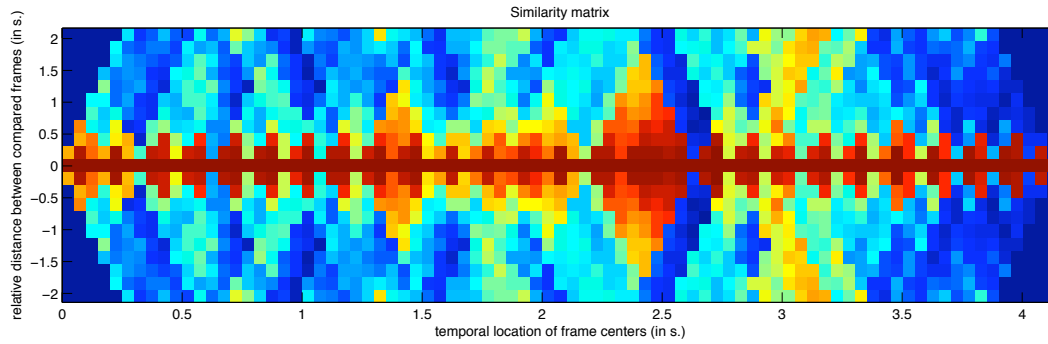
with a further normalization so that the values range from 0 to 1.



- `mirsimatrix(..., 'Horizontal')` rotates the matrix 45° in order to make the first diagonal horizontal:

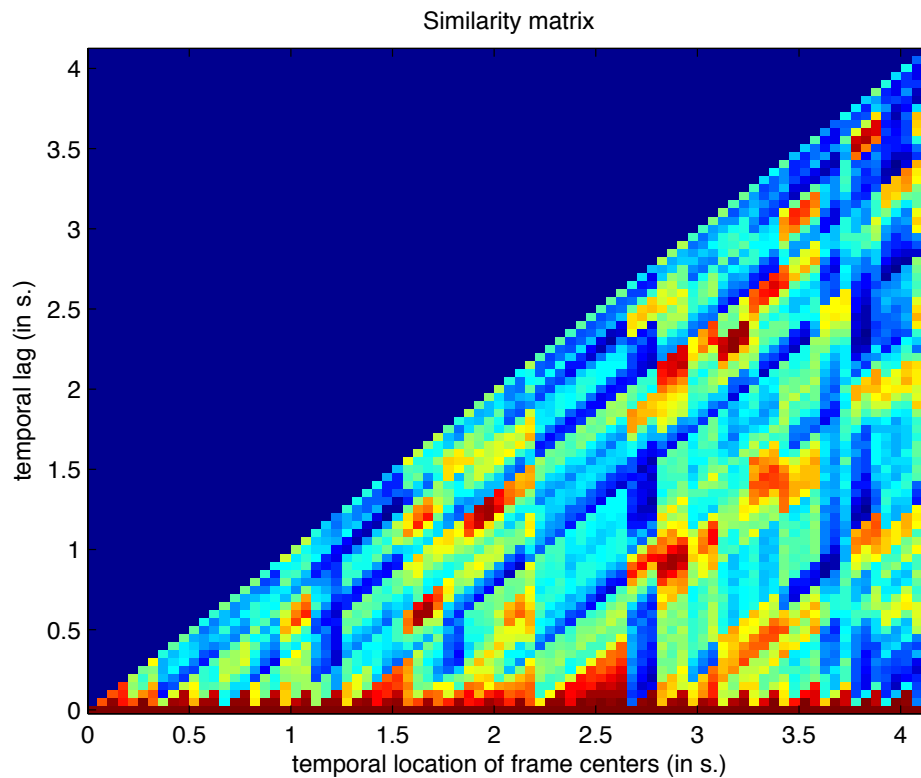


- `mirsimatrix(..., 'Horizontal', 'Width',  $w$ )` enables to restrict to a diagonal bandwidth only.

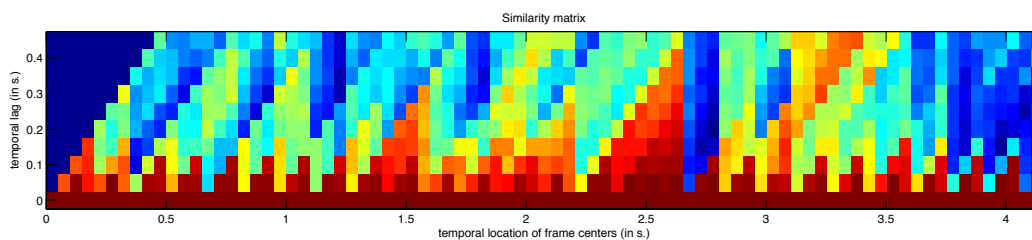


*mirsimatrix*(..., 'Width', 20, 'Horizontal')

- *mirsimatrix*(..., '**TimeLag**') computes the time-lag matrix out of the (non-rotated) initial time-time matrix:



- *mirsimatrix*(..., '**TimeLag**', '**Width**',  $w$ ) enables to restrict to a diagonal bandwidth only:



*mirsimatrix(..., 'Width', 20, 'TimeLag')*

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific field:

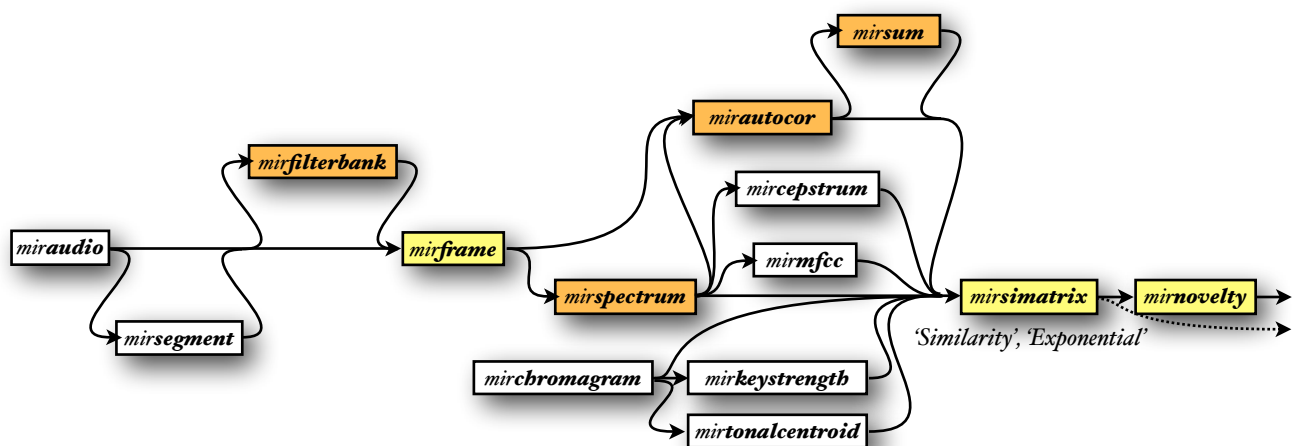
- ***'DiagWidth'***: the chosen value for the *'Width'* parameter.

## mirnovelty

### NOVELTY CURVE

One particular structural description related to the temporal succession of moments, each characterized by particular musical properties. The *novelty curve* represents the probability along time of the presence of transitions between successive states, indicated by peaks, as well as their relative importance, indicated by the peak heights.

### FLOWCHART INTERCONNECTIONS



Some parameters related to *mirsimatrix* are accessible in *mirnovelty*:

- **‘Distance’**,
- **‘Similarity’** (set here by default to *‘exponential’*),
- **‘Width’** (can also be called **‘KernelSize’** or **‘Kernel’**), with the default value here set to *Inf*. If no value is specified, our new multi-granular strategy for novelty computation (cf. below) is used. If on the contrary a value is specified, the kernel-based approach (cf. below) is chosen.
- **‘Horizontal’** and **‘TimeLag’** can be toggled on (both toggled off in the multi-granular strategy; in the kernel-based approach, the *‘TimeLag’* representation is chosen by default).

*mirnovelty* usually accepts either:

- *mirsimatrix* objects,
- *mirspectrum* frame-decomposed objects,
- *miraudio* objects (same as for *mirsimatrix*).
- **file name** or the **‘Folder’** keyword: same behavior than for *miraudio* objects;

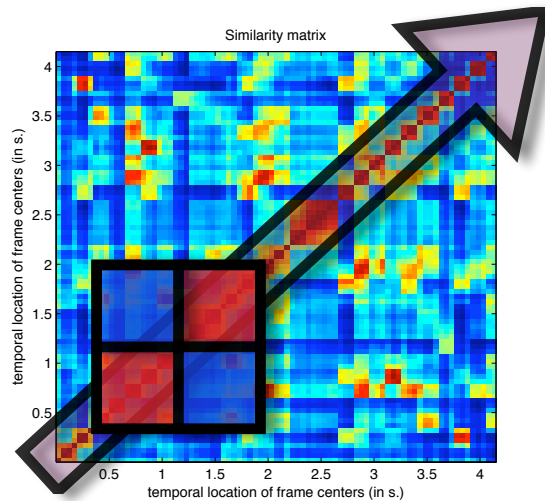
- ***mirautocor*** frame-decomposed objects;
- ***mircepstrum*** frame-decomposed objects;
- ***mirmfcc*** frame-decomposed objects;
- ***mirchromagram*** frame-decomposed objects;
- ***mirkeystrength*** frame-decomposed objects.

*mirnovelty* can return several outputs:

1. the novelty curve itself, and
2. the similarity matrix (***mirsimatrix***), horizontal.

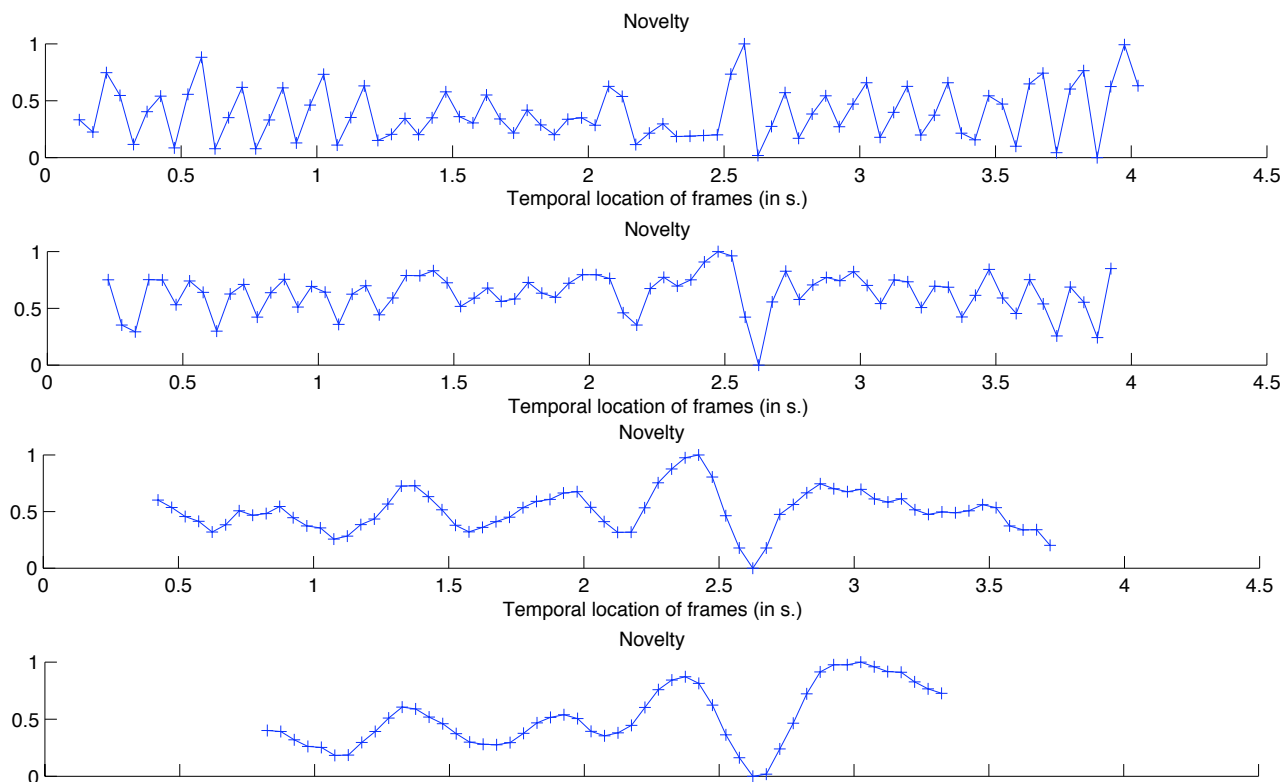
## 1. KERNEL - BASED APPROACH

Novelty is traditionally computed by comparing – through cross-correlation – local configurations along the diagonal of the similarity matrix with an ideal Gaussian checkerboard kernel (Foote & Cooper, 2003).



This method is dependent on the kernel size, focusing hence on a single level of granularity. The kernel size is expressed in number of samples in the similarity matrix. A kernel size of  $N$  means that it covers  $N$  rows and  $N$  columns in the similarity matrix.

Below is an example of novelty curve computed using increasing ‘Kernel’ size:



To use this kernel-based approach, a value for the ‘*Kernel*’ size parameter has to be specified. If no value is specified, the multi-granular approach, presented in the next paragraph, is chosen instead:

## 2 . M U L T I - G R A N U L A R   A P P R O A C H

Compared to the kernel-based approach, (Lartillot, Cereghetti, Eliard & Grandjean, 2013) introduces a simpler but more powerful and general method that automatically detects homogeneous segments of any size. For each instant in the piece, novelty is assessed by first determining the temporal scale of the preceding homogeneous part as well as the degree of contrast between that previous part and what just comes next. The idea is to estimate the temporal scale of the previous ending segment as well as the contrastive change before and after the ending of the segment. The novelty value is then represented as a combination of the temporal scale and the amount of contrast.

When the multi-granular approach is used for academic research, please cite the following publication:

Lartillot, O., Cereghetti, D., Eliard, K., Grandjean, D., “A simple, high-yield method for assessing structural novelty”, *International Conference on Music and Emotion*, Jyväskylä, 2013.



## POST-PROCESSING OPERATION

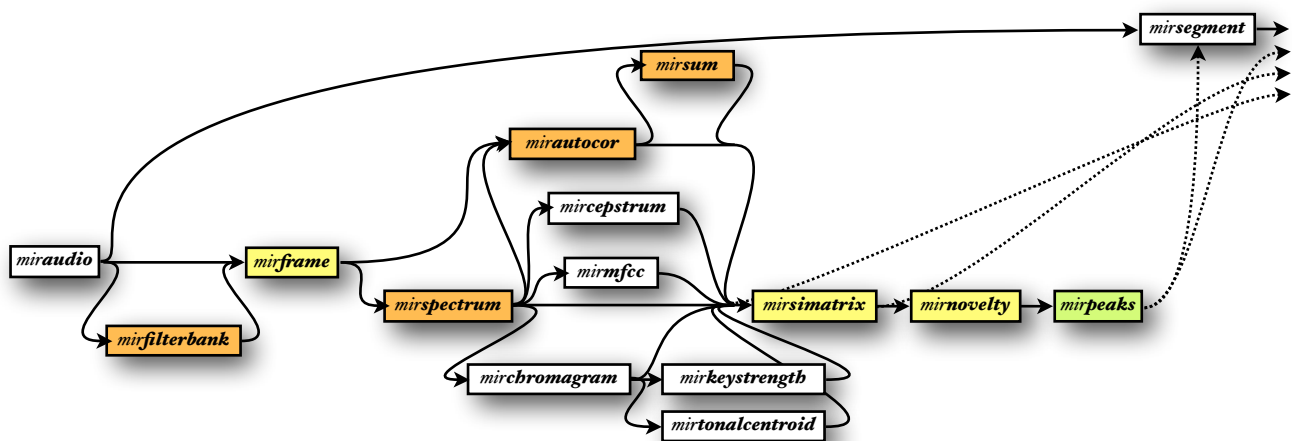
- *mirnovelty*(..., '**Normal**', *n*) toggles on/off the normalization of the novelty curve between the values 0 and 1. Toggled on by default.

## *mirsegment(..., 'Novelty')*

Peak detection applied to the novelty curve returns the temporal position of feature discontinuities that can be used for the actual segmentation of the audio sequence.

The 'Novelty' keyword is actually not necessary, as this strategy is chosen by default in *mirsegment*.

### FLOWCHART INTERCONNECTIONS



Some parameters related to *mirnovelty* are accessible in *mirsegment*: '**Distance**', '**Measure**' and '**KernelSize**'. Some parameters related to *mirpeaks* are accessible in *mirsegment*: '**Total**' (set by default to *Inf*) and '**Contrast**' (set to .1).

The choice of the feature used for the similarity matrix computation can be specified:

- *mirsegment(..., 'Spectrum')* will compute a *mirspectrum*, where some parameters can be specified: '*Min*', '*Max*', '*Normal*', '*Window*', '*Bark*', '*Mel*'. The default frame length is 50 ms and no overlapping.
- *mirsegment(..., 'MFCC')* will compute a *mirmfcc*, where the '*Rank*' parameter can be specified. Same default frame parameters than for '*Spectrum*'.
- *mirsegment(..., 'KeyStrength')* will compute a *mirkeystrength*. The default frame length is 500 ms and a hop factor of .2
- *mirsegment(..., 'Pitch')* ou *mirsegment(..., 'AucotorPitch')* will compute a *mirpitch* operation and use its second output, i.e., the *mirautocor*, for the computation of the similarity matrix. The default frame length is 50 ms and a hop factor of .01
- If no feature is specified, the default feature used in *mirsimatrix* will be chosen, i.e., the spectrum (*mirspectrum*).

The default frame parameters can be changed using the ‘*WinLength*’ option (in second) and the ‘*Hop*’ option (a value between 0 and 1).

*mirsegment* accepts uniquely as main input a ***miraudio*** objects not frame-decomposed, not channel decomposed, and not already segmented. Alternatively, **file name** or the ‘***Folder***’ key-word can be used as well.

*mirsegment*(..., ‘*Novelty*’) can return several outputs:

1. the segmented audio waveform itself,
2. the novelty curve (***mirnovelty***) after peak picking (***mirpeaks***),
3. the similarity matrix (***mirsimatrix***), and
4. the features entered into the ***mirsimatrix*** operator.

## 4.2. Statistics

### ***mirmean***

#### DESCRIPTION

*mirmean(f)* returns the mean along frames of the feature *f*.

If *f* is decomposed into segments, the results is the series of means in successive segments.

*f* can be a structure array composed of features. In this case, the output will be structured the same way.

### ***mirstd***

#### DESCRIPTION

*mirmean(f)* returns the standard deviation along frames of the feature *f*.

*f* can be a structure array composed of features. In this case, the output will be structured the same way.

### ***mirmedian***

#### DESCRIPTION

*mirmedian(f)* returns the median along frames of the feature *f*.

If *f* is decomposed into segments, the results is the series of medians in successive segments.

*f* can be a structure array composed of features. In this case, the output will be structured the same way.

## *mirstat*

### DESCRIPTION

*mirstat* can be applied to any object and will return its statistics in a structure array.

- If the object is frame-decomposed, the fields of the output are:
  - *Mean*: the average along frames;
  - *Std*: the standard deviation along frames;
  - *Slope*: the linear slope of the trend along frames, i.e. the derivative of the line that would best fit the curve. The slope  $S$  is computed in a normalised representation of the curve  $C$  – i.e., centered, with unit variance, and with a temporal scale reduced to a series  $T$  of values between  $-0.5$  and  $+0.5$  – as a solution in a least-square sense of the equation  $S*T = C$ ;
  - *PeriodFreq*: the frequency (in Hz.) of the maximal periodicity detected in the frame-by-frame evolution of the values, estimated through the computation of the autocorrelation sequence. The first descending slope starting from zero lag is removed from this analysis, as it is not related to the periodicity of the curve. If no periodicity is detected – i.e., if there is absolutely no peak in the autocorrelation sequence –, *NaN* is returned;
  - *PeriodAmp*: the normalized amplitude of that main periodicity, i.e., such that the autocorrelation at zero lag is identically 1.0.
  - *PeriodEntropy*: the Shannon entropy of the autocorrelation function (cf. *mirentropy*).
- If the object is not frame-decomposed, the data itself is returned directly in the single field *Mean*.

An additional field *FileNames* indicates the file name following the same order used when displaying each numerical result.

### N A N - F I L T E R

*mirstat* automatically discard any *NaN* value contained in the input data.

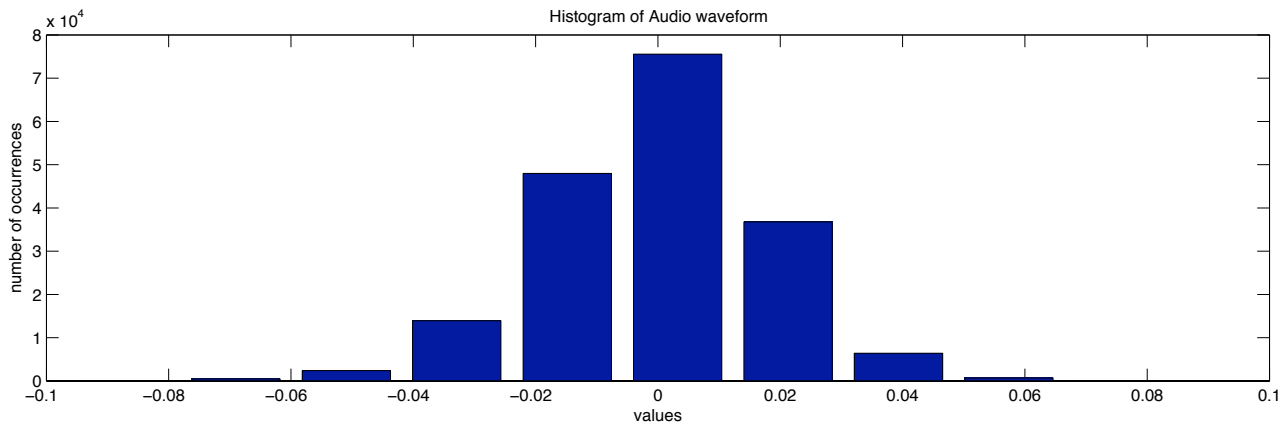
### M A N A G E M E N T   O F   S T R U C T U R E   A R R A Y S

If the input is already a structure array, the output will follows the same field decomposition (and all subfields as well) of the structure array, and will replace each final node with its corresponding *mirstat* structure array result.

## *mirbisto*

### DESCRIPTION

*mirbisto* can be applied to any object and will return its corresponding histogram. The data is binned into equally spaced containers.



### OPTIONS

- *mirbisto*(..., '**Number**', *n*) specifies the number of containers. Default value : *n* = 10.
- *mirbisto*(..., '**Ampli**') adds the amplitude of the elements, instead of simply counting them.

### ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- '**Weight**': the number of elements associated to each container (same as '*Data*'),
- '**Bins**': the range of value associated to each bin (same as '*Pos*'). A first layer of values (along the third matrix dimension) indicates the minimal value for each bin, a second layer indicates the maximal values.

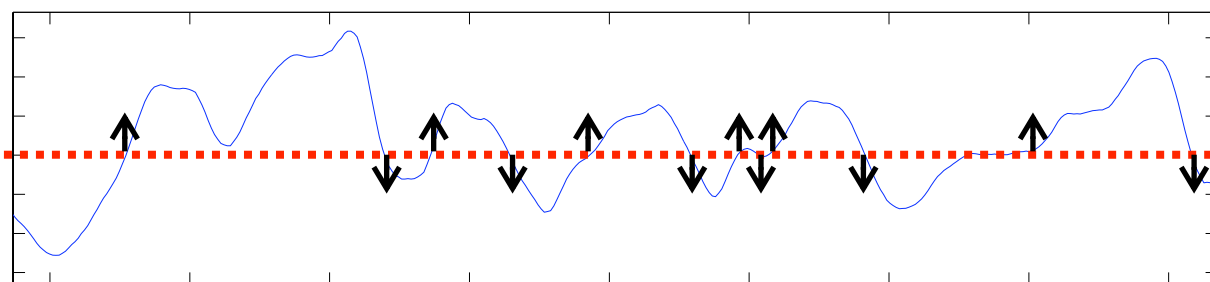
## *mirzerocross*

### DESCRIPTION

*mirzerocross* counts the number of times the signal crosses the X-axis (or, in other words, changes sign).

This function has already defined in the timbre section (§3.3): applied directly to audio waveform, *mirzerocross* is an indicator of noisiness.

But actually *mirzerocross* accepts any input data type.



### OPTIONS

- *mirzerocross*(..., '**Per**', *p*) precises the temporal reference for the rate computation. Possible values:
  - *p* = 'Second': number of sign-changes per second (Default).
  - *p* = 'Sample': number of sign-changes divided by the total number of samples. The 'Second' option returns a result equal to the one returned by the 'Sample' option multiplied by the sampling rate.
- *mirzerocross*(..., '**Dir**', *d*) precises the definition of sign change. Possible values:
  - *d* = 'One': number of sign-changes from negative to positive only (or, equivalently, from positive to negative only). (Default)
  - *d* = 'Both': number of sign-changes in both ways. The 'Both' option returns a result equal to twice the one returned by the 'One' option.

## *mircentroid*

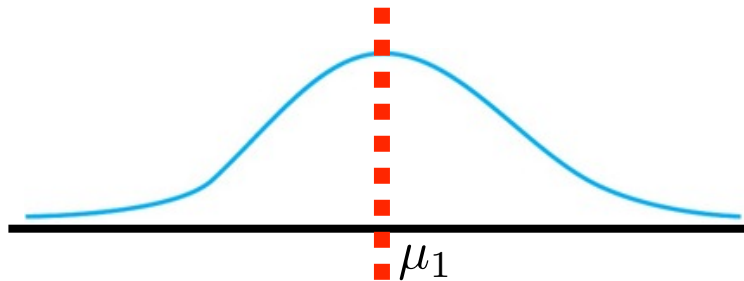
### DESCRIPTION

*mircentroid* returns the *centroid* of the data.

### EXPLANATION

An important and useful description of the shape of a distribution can be obtained through the use of its moments. The *first moment*, called the *mean*, is the geometric center (*centroid*) of the distribution and is a measure of central tendency for the random variable.

$$\mu_1 = \int x f(x) dx$$



### INPUTS

Any data can be used as input.

If the input is an audio waveform, a file name, or the ‘*Folder*’ keyword, the centroid is computed on the spectrum (spectral centroid). For quasi-silent frames, as specified by the ‘*MinRMS*’ option, *NaN* is returned.

If the input is a series of peak lobes produced by *mirpeaks*(..., ‘*Extract*’), the centroid will be computed for each of these lobes separately.

### OPTION

- When the input contains peaks (using *mirpeaks*), *mircentroid*(..., ‘**Peaks**’, *i*) will compute the centroids of the distribution of peaks. The argument *i* accepts two arguments:
  - *i* = ‘*NoInterpol*’: the centroid is computed on the non-interpolated peaks (default choice),
  - *i* = ‘*Interpol*’: the centroid is computed on the interpolated peaks (cf. ‘*Interpol*’ option in *mirpeaks*).



- *mircentroid*(..., '**MinRMS**', *m*) specifies the threshold *m*, as a value from 0 to 1, for the detection of quasi-silent frames, for which no value is given. For a given frame, if the RMS of the input (for spectral centroid, the input is each frame of the spectrogram, etc.) is below *m* times the highest RMS value over the frames, *NaN* is returned. default value: *m* = .005
- *mircentroid*(..., '**MaxEntropy**', *b*) specifies the threshold *b*, as a value from 0 to 1, for the detection of quasi-flat frames, for which no value is given. For a given frame, if the entropy of the input (for spectral centroid, the input is each frame of the spectrogram, etc.) is over *m*, *NaN* is returned. default value: *m* = .95

## *mirspread*

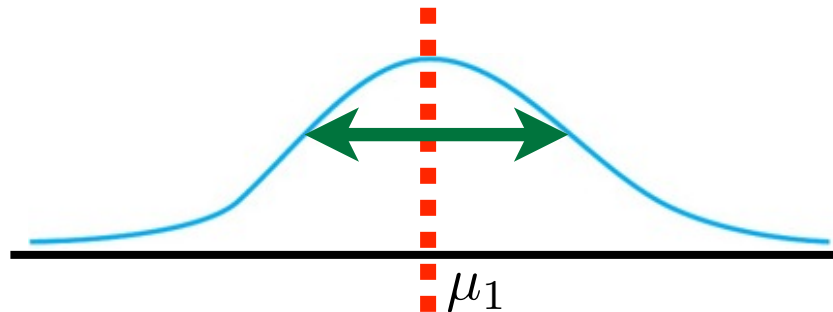
### DESCRIPTION

*mirspread* returns the *standard deviation* of the data.

### EXPLANATION

The *second central moment*, called the *variance*, is usually given the symbol sigma squared and is defined as:

$$\sigma^2 = \mu_2 = \int (x - \mu_1)^2 f(x) dx$$



Being the squared deviation of the random variable from its mean value, the variance is always positive and is a measure of the dispersion or spread of the distribution. The square root of the variance is called the *standard deviation*, and is more useful in describing the nature of the distribution since it has the same units as the random variable. (Koch)

### INPUTS

Any data can be used as input.

If the input is an audio waveform, a file name, or the ‘*Folder*’ keyword, the spread is computed on the spectrum (spectral spread).

If the input is a series of peak lobes produced by *mirpeaks*(..., ‘*Extract*’), the spread will be computed for each of these lobes separately.

### OPTION

- *mirspread*(..., ‘**MinRMS**’, *m*) specifies the threshold *m*, as a value from 0 to 1, for the detection of quasi-silent frames, for which no value is given. For a given frame, if the RMS of the input (for spectral spread, the input is each frame of the spectrogram, etc.) is below *m* times the highest RMS value over the frames, NaN is returned. default value: *m* = .01

## *mirskewness*

### DESCRIPTION

*mirskewness* returns the *coefficient of skewness* of the data.

### EXPLANATION

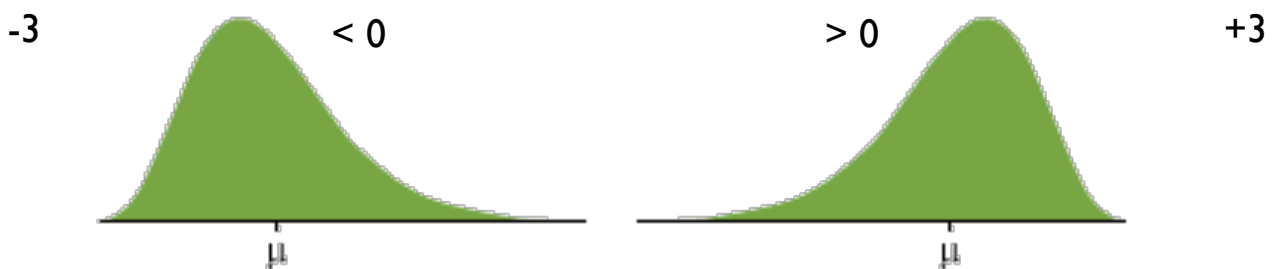
The *third central moment* is called the *skewness* and is a measure of the symmetry of the distribution. The skewness can have a positive value in which case the distribution is said to be positively skewed with a few values much larger than the mean and therefore a long tail to the right. A negatively skewed distribution has a longer tail to the left. A symmetrical distribution has a skewness of zero. (Koch)

$$\mu_3 = \int (x - \mu_1)^3 f(x) dx$$

The *coefficient of skewness* is the ratio of the skewness to the standard deviation raised to the third power.

$$\frac{\mu_3}{\sigma^3}$$

The coefficient of skewness has more convenient units than does the skewness and often ranges from -3.0 to 3.0 for data from natural systems. Again, a symmetrical distribution has a coefficient of skewness of zero. A positive coefficient of skewness often indicates that the distribution exhibits a concentration of mass toward the left and a long tail to the right whereas a negative value generally indicates the opposite. (Koch)



### INPUTS

Any data can be used as input.

If the input is an audio waveform, a file name, or the 'Folder' keyword, the skewness is computed on the spectrum (spectral skewness).

If the input is a series of peak lobes produced by *mirpeaks*(..., 'Extract'), the skewness will be computed for each of these lobes separately.

## *mirkurtosis*

### DESCRIPTION

*mirkurtosis* returns the (*excess*) *kurtosis*, of the data.

### EXPLANATION

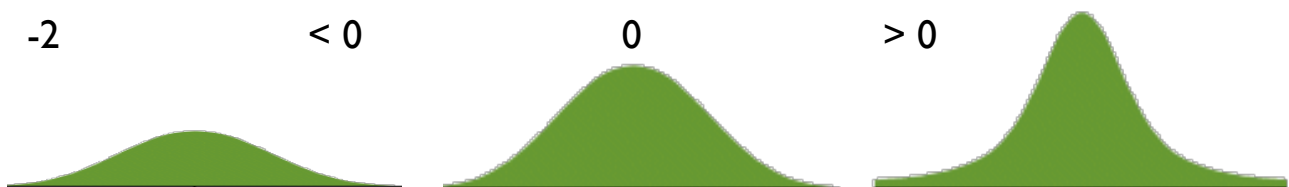
The *fourth standardized moment* is defined as,

$$\frac{\mu_4}{\sigma^4}$$

*Kurtosis* is more commonly defined as the fourth cumulant divided by the square of the variance of the probability distribution, equivalent to:

$$\frac{\mu_4}{\sigma^4} - 3$$

which is known as *excess kurtosis*. The "minus 3" at the end of this formula is often explained as a correction to make the kurtosis of the normal distribution equal to zero. Another reason can be seen by looking at the formula for the kurtosis of the sum of random variables. Because of the use of the cumulant, if  $\mathcal{Y}$  is the sum of  $n$  independent random variables, all with the same distribution as  $X$ , then  $Kurt[\mathcal{Y}] = Kurt[X]/n$ , while the formula would be more complicated if kurtosis were simply defined as *fourth standardized moment*. (Wikipedia)



### INPUTS

Any data can be used as input.

If the input is an audio waveform, a file name, or the 'Folder' keyword, the kurtosis is computed on the spectrum (spectral kurtosis).

If the input is a series of peak lobes produced by *mirpeaks*(..., 'Extract'), the kurtosis will be computed for each of these lobes separately.

## *mirflatness*

### DESCRIPTION

*mirflatness* returns the *flatness* of the data.

### EXPLANATION

The *flatness* indicates whether the distribution is smooth or spiky, and results from the simple ratio between the geometric mean and the arithmetic mean:

$$\frac{\sqrt[N]{\prod_{n=0}^{N-1} x(n)}}{\left(\frac{\sum_{n=0}^{N-1} x(n)}{N}\right)}$$

### INPUTS

Any data can be used as input.

If the input is an audio waveform, a file name, or the ‘*Folder*’ keyword, the flatness is computed on the spectrum (spectral flatness).

### OPTION

- *mirflatness*(..., ‘**MinRMS**’, *m*) specifies the threshold *m*, as a value from 0 to 1, for the detection of quasi-silent frames, for which no value is given. For a given frame, if the RMS of the input (for spectral flatness, the input is each frame of the spectrogram, etc.) is below *m* times the highest RMS value over the frames, NaN is returned. default value: *m* = .01

## *mirentropy*

### DESCRIPTION

*mirentropy* returns the relative Shannon (1948) entropy of the input. The Shannon entropy, used in information theory, is based on the following equation:

$$H(X) := - \sum_{i=1}^n p(x_i) \log_b p(x_i)$$

where  $b$  is the base of the logarithm.

In order to obtain a measure of entropy that is independent on the sequence length, *mirentropy* actually returns the relative entropy, computed as follows:

$$H(p) = -\text{sum}(p.*\log(p)) / \log(\text{length}(p));$$

Shannon entropy offers a general description of the input curve  $p$ , and indicates in particular whether it contains predominant peaks or not. Indeed, if the curve is extremely flat, corresponding to a situation of maximum uncertainty concerning the output of the random variable  $X$  of probability mass function  $p(x_i)$ , then the entropy is maximal. Reversely, if the curve displays only one very sharp peak, above a flat and low background, then the entropy is minimal, indicating a situation of minimum uncertainty as the output will be entirely governed by that peak.

The equation of Shannon entropy can only be applied to functions  $p(x_i)$  that follow the characteristics of a probability mass function: all the values must be non-negative and sum up to 1. Inputs of *mirentropy* are transformed in order to respect these constraints:

- Negative values are replaced by zeros (i.e., half-wave rectification).
- The remaining data is scaled such that it sums up to 1.

### INPUTS

Any data can be used as input.

If the input is an audio waveform, a file name, or the ‘*Folder*’ keyword, the entropy is computed on the spectrum (spectral entropy).

### OPTION

- *mirentropy*(..., ‘**Center**’) centers the input data before half-wave rectification.

- *mirentropy*(..., '**MinRMS**', *m*), when the input is an audio waveform, a spectrum or a cepstrum, specifies the threshold *m*, as a value from 0 to 1, for the detection of quasi-silent frames, for which no value is given. For a given frame, if the RMS of the input (for spectral entropy, the input is each frame of the spectrogram, etc.) is below *m* times the highest RMS value over the frames, *NaN* is returned. default value: *m* = .005



## *mirfeatures*

### DESCRIPTION

*mirfeatures* computes a large set of features, and returns them in a structure array. If the result is stored in a variable *f*, for instance, then the features are organized as follows:

- in a *dynamics* field,
  - *f.dynamics.rms{1}*: the frame-based RMS (*mirrms*), that can be accessed using the command:
- a *fluctuation* field, containing:
  - *f.fluctuation.peak{1}*: a fluctuation summary (*mirfluctuation*) with its highest peak (*mirpeak*),
  - *f.fluctuation.centroid{1}*: the centroid (*mircentroid*) of the fluctuation summary;
- in a *rhythm* field,
  - a *tempo* field, containing:
    - *f.rhythm.tempo{1}*: a frame-based tempo estimation (*mirtempo*),
    - *f.rhythm.tempo{2}*: the autocorrelation function used for the tempo estimation;
  - an *attack* field, containing:
    - a *time* field, with:
      - *f.rhythm.attack.time{1}*: the attack times (*mirattacktime*) of the events,
      - *f.rhythm.attack.time{2}*: the envelope curve used for the event detection (*mirevents*);
    - a *slope* field, with:
      - *f.rhythm.attack.slope{1}*: the attack slopes (*mirattackslope*) of the events;
- in a *timbre* field,
  - *f.timbre.zerocross{1}*: the frame-decomposed zero-crossing rate (*mirzerocross*),
  - *f.timbre.centroid{1}*: the frame-decomposed spectral centroid (*mircentroid*),
  - *f.timbre.brightness{1}*: the frame-decomposed brightness (*mirbrightness*),
  - *f.timbre.spread{1}*: the frame-decomposed spectral spread (*mirspread*),
  - *f.timbre.skewness{1}*: the frame-decomposed spectral skewness (*mirskewness*),
  - *f.timbre.kurtosis{1}*: the frame-decomposed spectral kurtosis (*mirkurtosis*),

- *f.timbre.rolloff95{1}*: the frame-decomposed roll-off (*mirrolloff*), using a 95 % threshold,
- *f.timbre.rolloff85{1}*: the frame-decomposed roll-off, using a 85 % threshold,
- *f.timbre.spectentropy{1}*: the frame-decomposed spectral entropy (*mirentropy*),
- *f.timbre.flatness{1}*: the frame-decomposed spectral flatness (*mirflatness*),
- a *roughness* field, containing:
  - *f.timbre.roughness{1}*: the frame-decomposed roughness (*mirroughness*),
  - *f.timbre.roughness{2}*: the spectrogram, containing the peaks used for the roughness estimation;
- an *irregularity* field, containing:
  - *f.timbre.irregularity{1}*: the frame-decomposed irregularity (*mirirregularity*),
  - *f.timbre.irregularity{2}*: the spectrogram, containing the peaks used for the irregularity estimation;
- an *inharmonic* field, containing:
  - *f.timbre.inharmonicity{1}*: the frame-decomposed inharmonicity (*mirinharmonicity*),
  - *f.timbre.inharmonicity{2}*: the spectrogram used for the inharmonicity estimation;
- *f.timbre.mfcc{1}*: the frame-decomposed MFCCs (*mirmfcc*),
- *f.timbre.dmfcc{1}*: the frame-decomposed delta-MFCCs,
- *f.timbre.ddmfcc{1}*: the frame-decomposed delta-delta-MFCCs,
- a *lowenergy* field, containing:
  - *f.timbre.lowenergy{1}*: the low energy rate (*mirlowenergy*),
  - *f.timbre.lowenergy{2}*: the RMS energy curve used for the low energy rate estimation;
- *f.spectralflux{1}*: the frame-decomposed spectral flux (*mirflux*);
- in a *pitch* field,
  - *f.pitch.salient{1}*: the frame-decomposed pitches (*mirpitch*),
  - a *chromagram* field, containing:
    - *f.pitch.chromagram.peak{1}*: an unwrapped chromagram (*mirchromagram*) and its highest peak,

- *f.pitch.chromagram.centroid{I}* field: the centroid of the chromagram;
- in a *tonal* field,
  - *f.tonal.keyclarity{I}*: the frame-decomposed key clarity (second output of *mirkey*),
  - *f.tonal.mode{I}*: the frame-decomposed mode (*mirmode*),
  - *f.tonal.hcdf{I}*: the frame-decomposed HCDF (*mirhcdf*).

## OPTION

- *mirfeatures(..., 'Stat')* returns the statistics of the features instead of the complete features themselves. In this way, the complete results are not accumulated in the RAM, preventing memory shortage problems.
- *mirfeatures(..., 'Segment', t)* segments the audio sequence at the temporal positions indicated in the array *t* (in s.), and analyzes each segment separately.

## *mirmap*

When *mirmap* is used for academic research, please cite the following publications:

Olivier Lartillot, Tuomas Eerola, Petri Toivainen, Jose Fornari, "Multi-feature modeling of pulse clarity: Design, validation, and optimization", *International Conference on Music Information Retrieval*, Philadelphia, 2008.

### DESCRIPTION

*mirmap(predictors\_file, ratings\_file)* performs a statistical mapping between ratings associated to a set of audio recordings, and a set of variables computed for the same set of audio recordings. It might be useful in particular when the set of predictors is large and diverse.

- The predictors were saved in one or several text file(s) whose name *predictors\_file* (or a cell array of file name) is given as first argument of the *mirmap* function. This text file is the result of the exportation (using *mirexport*) of a chosen analysis performed using *MIRtoolbox*, on the folder of audio files.
- The ratings are saved in a one or several text file(s) whose name *ratings\_file* (or a cell array of file name) is given as second argument of the *mirmap* function. This text file contains a simple column of values, with one value for each file analyzed in the folder of audio files under study, using the same ordering of files as for the predictors, i.e., a lexicographic order of file name.

The routine detects all features not sufficiently normal, based on a Lilliefors test (using Statistics Toolbox's *lillietest*) at the 5% significance level. Those non-normal features are then normalized using an optimization algorithm that automatically finds optimal Box-Cox transformations () of the data, ensuring that their distributions become sufficiently Gaussian, which is a prerequisite for correlation estimation. Features hence normalized that are still not considered as normal (this times using a more tolerant significance level of 1%) are then excluded from further analyses. They are listed in the Command Window under the heading "Excluded:".

From the selected features, only those whose correlation with the ratings is sufficiently statistically significant (with a p-value lower than .05) are selected.

The selected features are ordered from the most correlated to the least correlated ones. Features that are not sufficiently independent with respect to the better scoring ones (with a normalized crosscorrelation exceeding .6) are deleted as well.

### EXAMPLE

*f = mirfeatures('Folder')*

*mirexport('mypredictors.txt', f)*

*mirmap('mypredictors.txt', 'myratings.txt')*

## OUTPUT

*mirmap* returns a structure array, with the following fields:

- *normal*: the correlation between the normalized features and the ratings:
  - *cor*: the correlation value,
  - *pval*: the related p value;
- *significant*: the correlation between the statistically significant features and the ratings:
  - *cor*: the correlation value,
  - *pval*: the related p value,
  - *inter*: the cross-correlations with the other features,
  - *fields*: the feature name,
  - the resulting *alpha* and *lambda* values of the Box-Cox optimization;
- *best*: the correlation between the independent best features and the ratings:
  - *index*: the indices of chosen features among the significant ones,
  - *fields*: the feature name,
  - the related Box-Cox *alpha* and *lambda* values.

## TEXT DISPLAY

Various information are output in the Command Window:

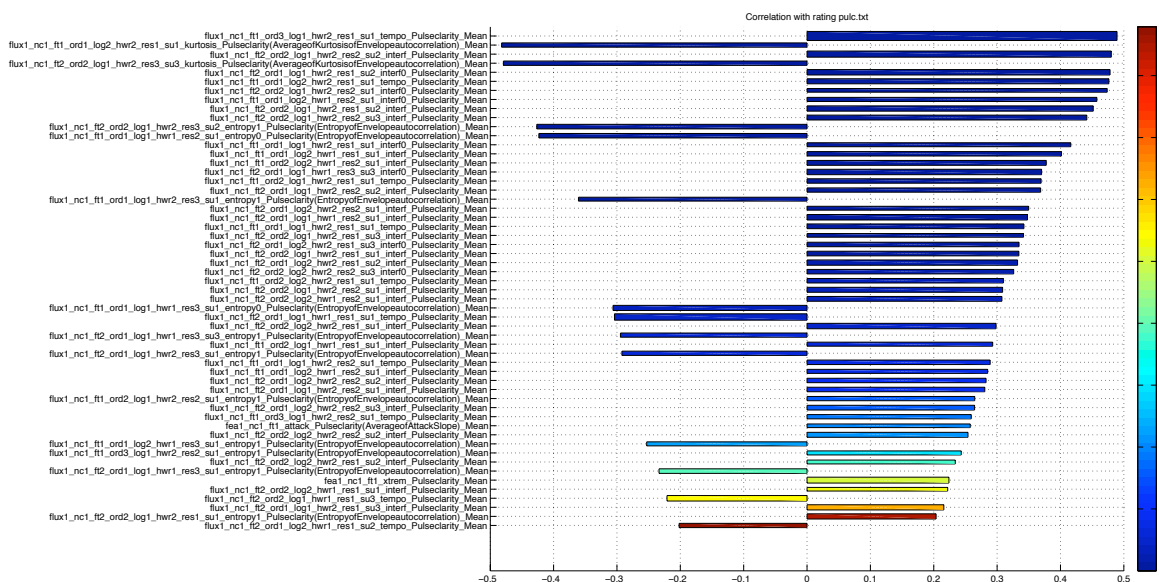
- the excluded features (cf. below)
- the finally selected features are then displayed, including:
  - the index among the significant features (*best.index*),
  - the feature name,
  - the correlation value with respect to the rating,
  - the worse cross-correlation with respect to better rating features.

- the result of a step-wise regression (using Statistics Toolbox *stepwisefit*), showing the successive steps, and the final results.
- the result of a normalized step-wise regression, this time centering and scaling each feature (using *stepwisefit* option ‘scale’).
- The list of features included in the final regression.

## GRAPHICAL OUTPUT

The correlation obtained by each of the best independent features are represented by a series of horizontal bars with best results at the top and worse results at the bottom. The abscissa axis represents the correlation values, that can be either positive or negative. The width of each bar indicates the independence of each feature with respect to the higher ones, and is computed as one minus the worse normalized cross-correlation. Finally, the color indicates the p-value, following the color scale given by the colorbar on the right of the figure: cold colors for statistically significant correlations, and warm colors for correlations of limited significance.

The figure below show a real-life example of useful application of the *mirmap* script. A large set of predictors were computed and stored in a highly structured array, composed of a hierarchical set of cells and structures corresponding to various methods for the predictor estimation.



### 4.3. Predictions

## *miremotion*

**The current version of *miremotion* is calibrated with version 1.3. Its use with more recent versions (1.3.1, ..., 1.4, ...) gives distorted results.**

When *miremotion* is used for academic research, please cite the following publication:

Tuomas Eerola, Olivier Lartillot, Petri Toivainen, "Prediction of Multidimensional Emotional Ratings in Music From Audio Using Multivariate Regression Models", *International Conference on Music Information Retrieval*, Kobe, 2009.

#### DESCRIPTION

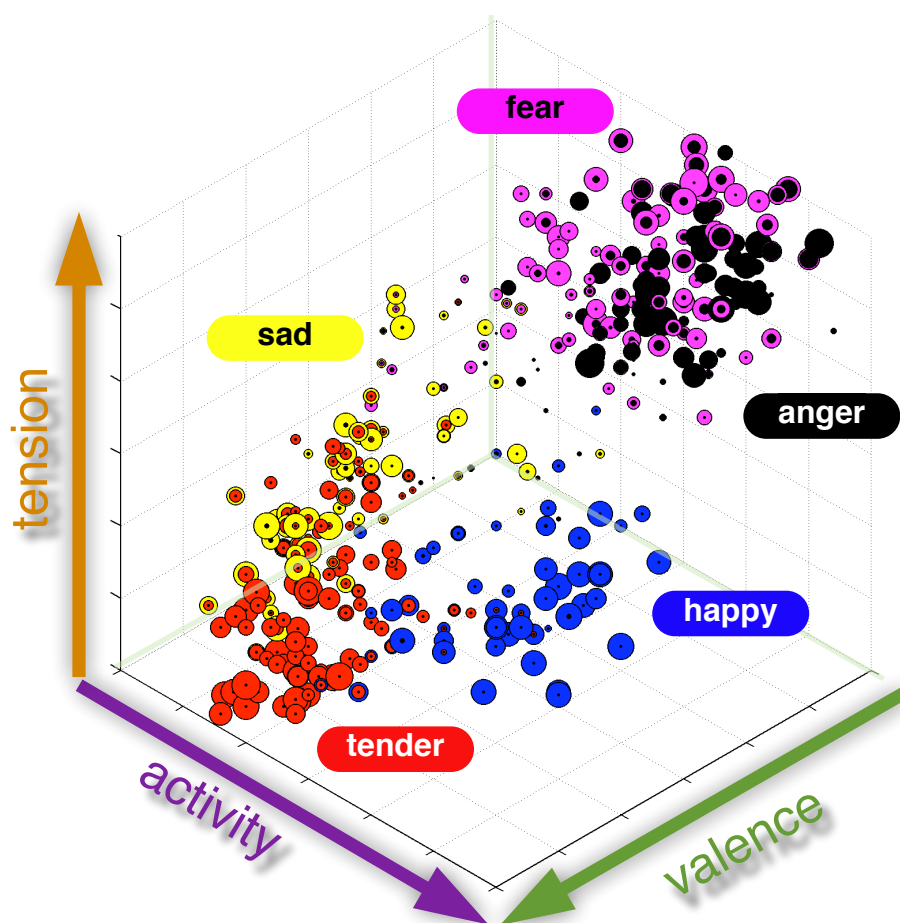
Emotions evoked in music is usually described along two opposite paradigms (cf. Eerola, Lartillot, Toivainen, 2009 for a literature review and a study of their interdependencies):

- Either as a decomposition into a few basic emotions. In *miremotion*, the 5 classes are happy, sad, tender, anger, fear.
- Either as a multi-dimensional space where the three dimensions, used in *miremotion*, are activity (or energetic arousal), valence (a pleasure-displeasure continuum) and tension (or tense arousal).

*miremotion* attempts to predict such description of emotion based on the analysis of the audio and musical contents of the recordings. Hence the output of *miremotion* corresponds to this underlying localization of emotional content within the 5 basic classes and within the 3 dimensions.

Each class or dimension is supposed to have values spanned in the interval [1,7]. Value 1 would correspond to very low value, and value 7 to very high value. This convention is used because the models are based on listeners' emotional ratings that were collected exactly the same way, as it corresponds to a classical Libert scale used in experimental psychology.

But even if the model was constructed using observation spanning in the interval [1,7], particular audio examples, not considered in the training, can extend beyond that range. So the interval [0,8] looks like a more probable range of value for these dimensions and concepts.



*Comparison between basic concepts of emotions (happy, sad, tender, anger, fear) and emotion dimensions (activity, valence, tension) (Eerola, Lartillot and Toiviainen, 2009).*

## PREDICTIVE MODELS

Following a systematic statistical methodology, a mapping has been found between each emotion concept or dimension and various sets of audio and musical features. The details of the procedure is given in (Eerola, Lartillot and Toiviainen, 2009).

The implemented models are based on multiple linear regression with 5 best predictors (MLR option in the paper). In this first version of *miremotion*, the Box-Cox transformations have been removed until the normalization values have been established with a large sample of music.

- The five factors contributing to the **activity** score are:
  - RMS averaged along frames ( $\beta = 0.6664$ )
  - Maximum value of summarized fluctuation ( $\beta = 0.6099$ )
  - Spectral centroid averaged along frames ( $\beta = 0.4486$ )



- Spectral spread averaged along frames ( $\beta = -0.4639$ )
- Entropy of the smoothed and collapsed spectrogram, averaged along frames ( $\beta = 0.7056$ )
- The five factors contributing to the **valence** score are:
  - Standard deviation of RMS along frames ( $\beta = -0.3161$ )
  - Maximum value of summarized fluctuation ( $\beta = 0.6099$ )
  - Key clarity (2nd output of *mirkey*) averaged along frames ( $\beta = 0.8802$ )
  - Mode averaged along frames ( $\beta = 0.4565$ )
  - Averaged spectral novelty ( $\beta = 0.4015$ )
- The five factors contributing to the **tension** score are:
  - Standard deviation of RMS along frames ( $\beta = 0.5382$ )
  - Maximum value of summarized fluctuation ( $\beta = -0.5406$ )
  - Key clarity (2nd output of *mirkey*) averaged along frames ( $\beta = -0.6808$ )
  - Averaged HCDF ( $\beta = 0.8629$ )
  - Averaged novelty from unwrapped chromagram ( $\beta = -0.5958$ )
- The five factors contributing to the **happy** score are:
  - Maximum value of summarized fluctuation ( $\beta = 0.7438$ )
  - Spectral spread averaged along frames ( $\beta = -0.3965$ )
  - Standard deviation of the position of the maximum of the unwrapped chromagram ( $\beta = 0.4047$ )
  - Key clarity (2nd output of *mirkey*) averaged along frames ( $\beta = 0.7780$ )
  - Mode averaged along frames ( $\beta = 0.6620$ )
- The five factors contributing to the **sad** score are:
  - Spectral spread averaged along frames ( $\beta = 0.4324$ )
  - Standard deviation of the position of the maximum of the unwrapped chromagram ( $\beta = -0.3137$ )
  - Mode averaged along frames ( $\beta = -0.5201$ )

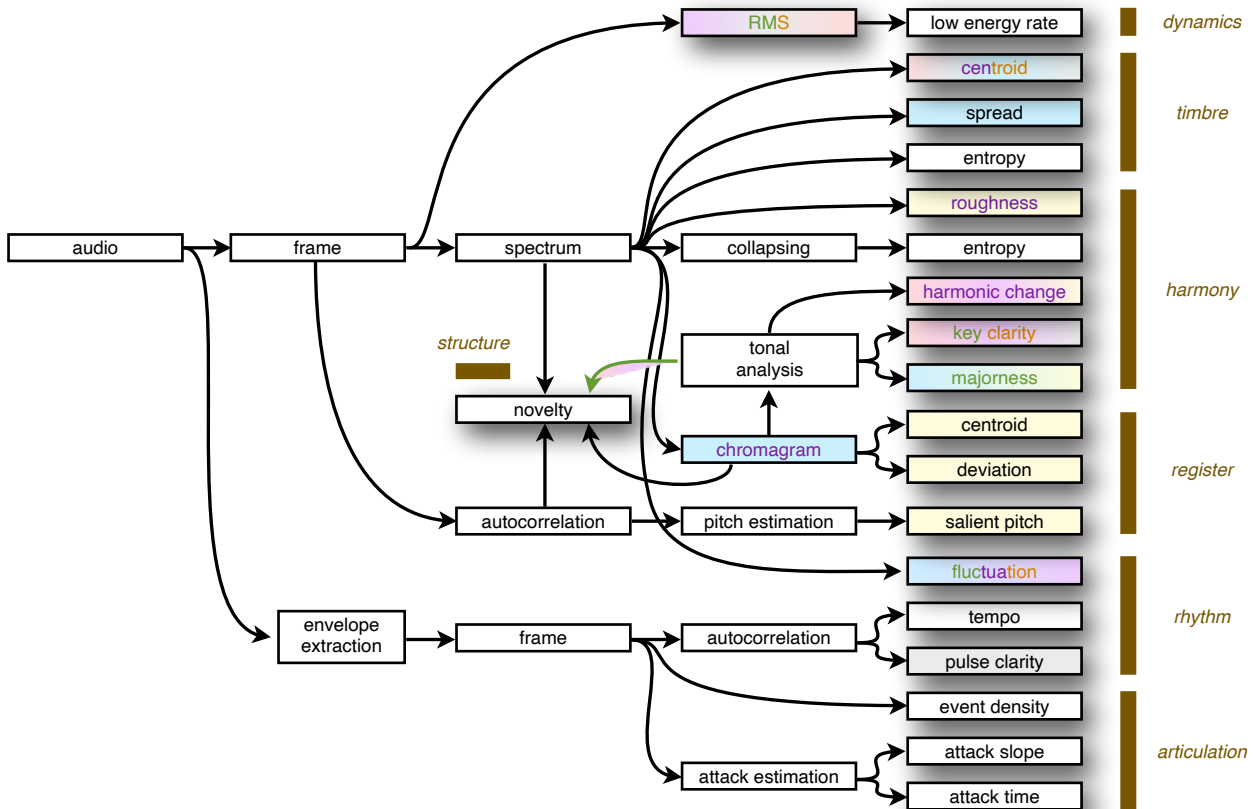
- Averaged HCDF ( $\beta = -0.6017$ )
- Averaged novelty from wrapped chromagram ( $\beta = 0.4493$ )
- The five factors contributing to the **tender** score are:
  - Spectral centroid averaged along frames ( $\beta = -0.2709$ )
  - Standard deviation of roughness ( $\beta = -0.4904$ )
  - Key clarity (2nd output of *mirkey*) averaged along frames ( $\beta = 0.5192$ )
  - Averaged HCDF ( $\beta = -0.3995$ )
  - Averaged spectral novelty ( $\beta = 0.3391$ )
- The five factors contributing to the **anger** score are:
  - Roughness averaged along frames ( $\beta = 0.5517$ )
  - Key clarity (2nd output of *mirkey*) averaged along frames ( $\beta = -0.5802$ )
  - Entropy of the smoothed and collapsed spectrogram, averaged along frames ( $\beta = 0.2821$ )
  - Averaged novelty from unwrapped chromagram ( $\beta = -0.2971$ )
- The five factors contributing to the **fear** score are:
  - Standard deviation of RMS along frames ( $\beta = 0.4069$ )
  - Averaged attack time ( $\beta = -0.6388$ )
  - Maximum value of summarized fluctuation ( $\beta = -0.2538$ )
  - Key clarity (2nd output of *mirkey*) averaged along frames ( $\beta = -0.9860$ )
  - Mode averaged along frames ( $\beta = -0.3144$ )

For later version of *MIRtoolbox*, we plan to revise the coefficients in order to

(a) force the output range between 0 - 1 and

(b) ground them on alternative models and materials (training sets).

A classifier of emotion concepts based also on audio features is also under development.



*Audio and musical features founding the prediction of emotion. The colors associated to each feature corresponds to the correlating emotions, following the same color code that in the previous figure (Eerola, Lartillot and Toivainen, 2009).*

## OPTIONAL ARGUMENTS

- *miremotion*(..., '**Frame**',  $l$ ,  $b$ ) predicts the emotional content for each successive frame of length  $l$ , and with a hop factor of  $b$ . By default,  $l = 1$  second and no overlapping.
- *miremotion*(..., '**Dimensions**',  $n$ ) indicates the number of emotion dimensions to output (0, 2 or 3). By default,  $n = 3$ . Alternatively, these dimensions can be listed:
  - *miremotion*(..., '**Activity**')
    - *miremotion*(..., '**Valence**')
      - *miremotion*(..., '**Tension**')

A value  $n = 2$  will output Activity and Valence only.

- *miremotion*(..., '**Arousal**') corresponds to *miremotion*(..., '**Activity**', '**Tension**').

- *miremotion(..., 'Concepts')* outputs all the emotion concepts (chosen by default). Alternatively, these concepts can be listed:
  - *miremotion(..., 'Happy')*
  - *miremotion(..., 'Sad')*
  - *miremotion(..., 'Tender')*
  - *miremotion(..., 'Anger')*
  - *miremotion(..., 'Fear')*
- *miremotion(..., 'Concepts', 0)* excludes all emotion concepts in the results.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **'Dim'**: the list of emotion dimensions taken into consideration,
- **'DimData'**: the score associated with each of these emotion dimensions,
- **'Class'**: the list of emotion classes taken into consideration,
- **'ClassData'**: the score associated with each of these emotion classes,
- **'ActivityFactors'**: the value associated with each of the factors contributing to the activity score (in the same order as presented above),
- **'ValenceFactors'**: the value associated with each of the factors contributing to the valence score (in the same order as presented above),
- **'TensionFactors'**: the value associated with each of the factors contributing to the tension score (in the same order as presented above),
- **'HappyFactors'**: the value associated with each of the factors contributing to the happy score (in the same order as presented above),
- **'SadFactors'**: the value associated with each of the factors contributing to the sad score (in the same order as presented above),
- **'TenderFactors'**: the value associated with each of the factors contributing to the tender score (in the same order as presented above),
- **'AngerFactors'**: the value associated with each of the factors contributing to the anger score (in the same order as presented above),

- ***FearFactors***: the value associated with each of the factors contributing to the fear score (in the same order as presented above).

## *mirclassify*

### DESCRIPTION

*mirclassify(features)* classifies along the analytic feature(s) *features* using cross-validation. *Features* can be either one feature, or a cell array of features. All features should be computed using ‘*Folders*’ on a same folder of folders audio files: the main folder should be composed of a certain number of sub-folders of same size, corresponding to the different classes. The name of each sub-folder will be used as the name of the corresponding class.

The classification is carried out using leave-one-out cross-validation.

*misclassify* requires the *Netlab* toolbox.

You can also integrate your own arrays of numbers computed outside *MIRtoolbox* as part of the features. These arrays should be given as matrices where each successive column is the analysis of each successive file.

### EXAMPLE

*mirclassify(mirmfcc('Folders'))*

*mirclassify({mirmfcc('Folders'), mircentroid('Folders')})*

### OPTIONAL ARGUMENTS

- *mirclassify(..., 'Nearest')* uses the minimum distance strategy. (by default)
- *mirclassify(..., 'Nearest', k)* uses the *k*-nearest-neighbour strategy. Default value: *k* = 1, corresponding to the minimum distance strategy.
- *mirclassify(..., 'GMM', ng)* uses a gaussian mixture model. Each class is modeled by at most *ng* gaussians. Default value: *ng* = 1.
  - Additionnally, the type of mixture model can be specified, using the set of value proposed in *Netlab*'s *gmm* function: i.e., '*spherical*', '*diag*', '*full*' (default value) and '*ppca*'. (cf. help *gmm*)

### ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- '**Classes**': the class associated to each test sample (same as '*Data*'),
- '**Correct**': the correctness of the classification, from 0 to 1, for each cross-validation fold.

## *mircluster*

### DESCRIPTION

- *mircluster(a, f, ...)*, where *a* is one or several audio files *already segmented* (using *mirsegment*), and *f* is an analysis of the same data, clusters the segments in the audio sequence(s) contained in the audio object *a*, along the analytic feature(s) *f*, using the *k*-means strategy.
  - The analytic feature(s) *f* should *not* be frame decomposed. Frame-decomposed data should first be summarized, using for instance *mirmean* or *mirstd*.
  - Multiple analytic features have to be grouped into one array of cells.
- *mircluster(d, ...)*, where *d* is any *frame-decomposed* feature computed using *MIRtoolbox*, clusters the segments in the audio sequence(s) contained in the audio object *a*, along the analytic feature(s) *f*, using the *k*-means strategy. Multiple analytic features have to be grouped into one array of cells.

*mircluster* simply calls the *kmeans\_clusters* function from the *SOM toolbox* (that is part of *MIRtoolbox* distribution) with the data contained in the *MIRtoolbox* objects provided as input. This is the only use of *SOM toolbox* in *MIRtoolbox*.

### EXAMPLE

- Clustering of audio segments:

*sg = mirsegment(a);*

*mircluster(sg, mirmfcc(sg))*

*mircluster(sg, {mirmfcc(sg), mircentroid(sg)})*

- Clustering of frames:

*cc = mirmfcc(a, 'Frame');*

*mircluster(cc)*

### OPTIONAL ARGUMENT:

- *mircluster(..., n)* indicates the maximal number of clusters. Default value: *n* = 2.
- *mircluster(..., 'Runs', r)* indicates the maximal number of runs. Default value: *r* = 5.

## ACCESSIBLE OUTPUT

cf. §5.2 for an explanation of the use of the *get* method. Specific fields:

- **‘*Clusters*’**: a structured descriptions of the clustered reduction (centroids, reduction of the initial representation as a sequence of clusters, etc.),



## 4.4. Similarity and Retrieval

Operators have been added for the comparison between audio files and for query by examples.

### ***mirdist***

#### DESCRIPTION

*mirdist(x,y)* evaluates the distance between audio files along a particular representation specified by the user, corresponding to audio or musical features computed using *MIRtoolbox*. The input variables *x* and *y* are the obtained representations computed for the chosen audio files.

*x* should correspond to the representation of one particular file, whereas *y* can correspond to the representation of either one particular file:

$$x = \text{mirmfcc}(\text{'file1.mp3'});$$
$$y = \text{mirmfcc}(\text{'file2.mp3'});$$
$$\text{mirdist}(x,y)$$

or a folder of files:

$$x = \text{mirmfcc}(\text{'file1.mp3'});$$
$$y = \text{mirmfcc}(\text{'Folder'});$$
$$\text{mirdist}(x,y)$$

- If *x* and *y* are not decomposed into frames,

*mirdist(..., metric)* specifies the distance method *metric*, following the same list as used in Matlab *pdist* command (cf. *help pdist*). Default value: *metric* = 'Cosine'.

- If *x* and *y* are composed of clustered frames (using *mircluster*), the cluster signatures are compared using Earth Mover Distance (Logan and Salomon, 2001).
- If *x* and *y* contains peaks, the vectors representing the peak distributions are compared using Euclidean distance (used with *mirnovelty* in Jacobson, 2006).

## *mirquery*

### DESCRIPTION

*mirquery*( $q, b$ ), where

- $q$  is the analysis of one audio file and
- $b$  is the analysis of a folder of audio files,

according to the same *MIRtoolbox* feature, returns the name of the audio files in the database  $b$  in an increasing distance to  $q$  with respect to the chosen feature, where the distance is computed using *mirdist*.

*mirquery*( $d$ ), where  $d$  is the distance between one audio file and a folder of audio file, according to a *MIRtoolbox* feature, returns the name of the audio files in an increasing distance  $d$ .

### OPTIONAL ARGUMENTS

- *mirquery*(..., '**Best**',  $n$ ) returns the name of the  $n$  closest audio files.
- *mirquery*(..., '**Distance**',  $metric$ ) specifies the distance method  $metric$  to use (cf. *mirdist*). Default value:  $metric = 'Cosine'$ .

The successive results can then be played using *mirplay*.

## 4.5. Exportation

### *mirgetdata*

#### DESCRIPTION

*mirgetdata* return the data contained in the input in a structure that can be used for further computation outside *MIRtoolbox*.

#### OUTPUT FORMAT

- If the input is based on one non-segmented audio sequence, the result is returned as a matrix. The columns of the matrix usually correspond to the successive frames of the audio signal. The third dimension of the matrix corresponds to the different channels of a filterbank.
- If the input corresponds to a set of audio sequences, and if each sequence has same number of frames, the corresponding resulting matrices are concatenated columnwise one after the other. If the number of rows of the elementary matrices varies, the missing values are replaced by NaN in the final matrix. On the contrary, if the number of columns (i.e., frames) differs, then the result remains a cell array of matrices.
- Idem if the input corresponds to one or several segmented audio sequence(s).

#### INPUT AND OUTPUT

- If the input is a key strength curve, the fourth dimension distinguishes between major and minor keys. i.e.:
  - $d(:, :, :, 1)$  is the keystrength for the major keys, and
  - $d(:, :, :, 2)$  is the keystrength for the minor keys.
- If the input is a key estimation, two output are returned:
  - a. the keys (from 1 to 12) and
  - b. the modes (1 for major, 2 for minor).
- If the input is the result of a peak detection, two output are returned:
  - a. the position of the peaks and
  - b. the value corresponding to these peaks, in the units predefined for this data.

What is returned by *mirgetdata* is no more a *MIRtoolbox* object, but just a *Matlab* array (or matrix). So of course, if you display such array using *Matlab plot* function, what you get is a curve without any proper X-axis annotation. (So by default the X-axis simply indicates the sample indices).

When using *mirgetdata* with *'Folder'* or *'Folders'*, the concatenation of data across audio files can be sometimes difficult to parse. In that case, it might be easier to use directly *get(..., 'Data')* (cf. §5.4).

## *mirexport*

*mirexport*(filename, ...) exports statistical information related to diverse data into a text file called filename.

*mirexport*(**'Workspace'**, ...) instead directly output the statistical information in a structure array saved in the Matlab workspace. This structure contains three fields:

- *filenames*: the name of the original audio files,
- *types*: the name of the features,
- *data*: the data.

The exported data should be related to the same initial audio file or the same ordered set of audio files.

The data listed after the first arguments can be:

- any feature computed in *MIRtoolbox*. What is exported is the statistical description of the feature (using the *mirstat* function).

➡ If the keyword *mirexport*(..., **'Raw'**) is added, what is exported is the complete, raw, feature data instead. Feature data are stored in columns, and arrays shorter than others are completed with NaN values in order to obtain a matrix.

- any structure array of such features. Such as the output of the *mirstat* function.
- any cell array of such features.
- the name of a text file. The text file is imported with the Matlab *importdata* command. Each line of the file should contain a fixed number of data delimited by tabulations. The first line, or 'header', indicates the name of each of these columns.

The file format of the output can be either:

- a text file. It follows the same text file representation as for the input text files. The first column of the matrix indicates the name of the audio files. The text file can be opened in Matlab, or in a spreadsheet program, such as Microsoft Excel, where the data matrix can be automatically reconstructed.
- an attribute-relation file. It follows the ARFF standard, used in particular in the WEKA data mining environment.

If the audio files were initially labelled using the *'Label'* option in *miraudio*, then these labels are automatically indicated as classes in the ARFF output file.



## 5. ADVANCED USE OF *MIRTOOLBOX*

### 5.1. *Parallel processing*

*mirparallel* does not work any more for Matlab 2013b and more recent. If you would like *MIRtoolbox* to be developed for parallel processing or for other purposes, and if you have any funding to suggest, please contact us.

When ‘*Folder*’ or ‘*Folders*’ is used, several audio files can be analysed in parallel using several parallel *Matlab* sessions running on the different processors and/or processor cores of your computer.

This requires MathWorks’ *Parallel Computing Toolbox*.

**M I R P A R A L L E L**

*mirparallel*(1) toggles on parallel processing.

*mirparallel*(*N*) with  $N > 1$  specifies the number of pools as *N*. If on the contrary *N* is equal to 1 or to Inf, the default number of pools specified by the default cluster profile is used instead.

*mirparallel*(0) toggles back off parallel processing.

*mirparallel* without argument returns the current setting, either 0, 1, or an integer *N* higher than 1.

## 5.2. *Progressive storage of stats*

When ‘*Folder*’ or ‘*Folders*’ is used, the statistical result of each audio file can be progressively added to a text file called *mirtemporary.txt*.

### M I R T E M P O R A R Y

*mirtemporary(1)* toggles on the progressive storage of stat results.

*mirtemporary(0)* toggles back off the progressive storage of stat results.

*mirtemporary* without argument returns the current setting, either 0 or 1.



### 5.3. *Interface preferences*

#### M I R V E R B O S E

*mirverbose(o)* toggles off the display by *MIRtoolbox* of minor informations in the Matlab Command Window (such as "*Computing mirfunction ...*").

*mirverbose(1)* toggles back on the display of such information.

*mirverbose* without argument returns the current setting, either 0 or 1.

#### M I R W A I T B A R

*mirwaitbar(o)* toggles off the display by *MIRtoolbox* of waitbar windows.

*mirwaitbar(1)* toggles back on the display of these waitbar windows.

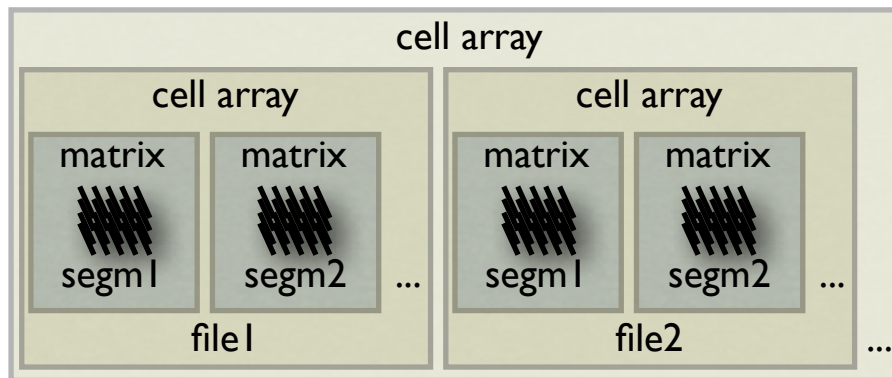
*mirwaitbar* without argument returns the current setting, either 0 or 1.

## 5.4. *get*

*get* returns fields of *MIRtoolbox* objects.

General fields used by most objects are the following:

- **'Data'**: the computed values, stored in a cell array, where each cell, corresponding to one audio file, is itself a cell array, where each cell, corresponding to one segment, is a matrix, where columns indicate successive frames, and where multiple channels are represented in the 3rd dimension.



- **'Pos'**: the X-axis abscissae related to the Y-axis ordinates given by 'Data', and stored in the same way as for 'Data'.
- **'Unit'**: the name of the unit in which the values are represented,
- **'Sampling'**: the sampling rate of the data, stored in a array of numbers, one number for each audio file.
- **'Length'**: the duration of the input audio signal, in number of samples.
- **'NBits'**: the resolutions in number of bits of the initial audio file, stored in the same ways as for 'Sampling'.
- **'Name'**: the name of the audio files, stored in a cell array, where each cell indicates the name of one file.
- **'Title'**: the name of the data type,
- **'Channels'**: the number and indexing of the channels decomposing the initial signal, stored in a cell array, where each cell contains the channel indexes for one audio file.

The list of more specific fields are indicated in the "*Accessible Output*" paragraph concluding the description of each feature.

## 5.5. Memory management

There are important things to know in order to take benefit of the memory management mechanism offered by *MIRtoolbox*.

### MIRCHUNKLIM

*mirchunklim* without argument returns the current chunk size (cf. below).

*mirchunklim*(*n*), where *n* is an integer, sets the chunk size to *n*. For instance, to increase 10 times the chunk size, you can write:

$$\textit{mirchunklim}(\textit{mirchunklim} * 10)$$

*mirchunklim*(*Inf*) toggles off the chunk decomposition.

### AUTOMATED MEMORY OPTIMIZATION

When a *MIRtoolbox* operator is called with ‘Folder’ (or ‘Folders’) as main argument, the command is applied to each audio file in the current directory one file after the other. The previous audio file is cleared from the memory when the next audio file is considered. The results for all audio file are combined into the main memory and returned. For that reason, ‘Folder’ (and ‘Folders’) keywords should be applied to the feature of highest possible level. For instance, when evaluating the key of a batch of file,

$$o = \textit{mirkey}(\textit{Folder})$$

each audio file is processed one after the other, and only the key result of each file is stored and grouped with those of the other files. This does not cause any memory problem.

On the contrary, the following set of commands should be avoided:

$$a = \textit{miraudio}(\textit{Folder})$$
$$o = \textit{mirkey}(a)$$

Why is this code problematic? Because after executing the first line, *all* the waveform of *all* the audio files are supposed to be stored together in the memory, which might sometimes leads to memory overflow problems.

Similarly, when a *MIRtoolbox* operator is called with used for a particularly long audio file, the audio file is automatically decomposed into little chunks that can fit into the memory. The

command is applied to each successive chunk one after the other. The previous chunk is cleared from the memory when the next chunk is considered. The results for all chunks are combined into the main memory and returned. For that reason, audio files should be applied to the feature of highest possible level. For instance, when evaluating the key evolution of long file,

$$o = \text{mirkey}(\text{'myhugefile'}, \text{'Frame'})$$

each successive chunk is processed one after the other, and only the key result of each chunk is stored and concatenated with those of the other chunks. This does not cause any memory problem.

On the contrary, the following set of commands should be avoided:

$$a = \text{miraudio}(\text{'myhugefile'})$$

$$o = \text{mirkey}(a, \text{'Frame'})$$

Why is this code problematic? Because after executing the first line, the *complete* waveform of the audio files is supposed to be stored entirely in the memory, which might sometimes leads to memory overflow problems.

For those reasons, as mentioned previously, when possible avoid writing a succession of operators if you can write in one line more efficiently.

Hopefully, there is a way to avoid such consideration by using the '*Design*' keyword, as will be explained below.

## FLOWCHART DESIGN

If you write this command:

$$a = \text{miraudio}(\text{'myhugefile'})$$

$$s = \text{mirspectrum}(a, \text{'Frame'})$$

$$c = \text{mircentroid}(s)$$

there may be memory problem when computing the spectrogram. In order to benefit from *MIRtoolbox* memory management mechanisms, you should write instead:

$$a = \text{miraudio}(\textbf{'Design'})$$

```
s = mirspectrum(a, 'Frame')
```

```
c = mircentroid(s)
```

```
mireval(c, 'myhugefile')
```

## COMPLEX FLOWCHART DESIGN

But if you want to get several output in your flowchart, for instance:

```
s = mirspectrum('mysong');
```

```
cent = mircentroid(s);
```

```
ceps = mircepstrum(s);
```

You should use a *mirstruct* object, which store all the outputs of your flowchart into fields and write as follows:

```
myflow = mirstruct
```

All the temporary data used in your flowchart should be stored into a *tmp* field. Hence when a spectrogram is used for several final features, it should be a temporary variable:

```
myflow = mirstruct
```

```
myflow.tmp.s = mirspectrum('Design', 'Frame');
```

```
myflow.cent = mircentroid(myflow.tmp.s);
```

```
myflow.ceps = mircepstrum(myflow.tmp.s);
```

```
output = mireval(myflow, 'myhugefile');
```

Similarly, when an event detection curve is used for several final features, it should be a temporary variable:

```
myflow = mirstruct;
```

```
myflow.tmp.o = mirevents('Design');
```

```
myflow.at = mirattacktime(myflow.tmp.o);
```

```
myflow.as = mirattackslope(myflow.tmp.o);
```

```
myflow = mirstat(myflow);
```

```
output = mireval(myflow, 'myhugefile');
```

Please note also that in the current version of *MIRtoolbox*, it is not possible to identify directly an output variable with a temporary variable, such as:

```
r.rms = r.dynamics.tmp.rms;    % Does not work yet.
```

As a temporary solution, you can call the *MIRtoolbox* operator once again, such as:

```
r.rms = mirrms(r.dynamics.tmp.rms);    % This is OK.
```

Note also that, as can be seen in the previous example, you can keep only the statistics of the results instead of the complete results by adding the line:

```
myflow = mirstat(myflow);
```

## STRUCTURED COMPLEX FLOWCHART DESIGN

You can decompose your flowchart into several fields, using the standard structure array used in *Matlab*, for instance associating each field with one particular musical dimension (such as 'dynamics', etc.). You can assign a temporary variable within one particular field, by declaring that field as a *mistruct* object:

```
r.dynamics = mistruct;
```

```
r.dynamics.tmp.rms = mirrms(a, 'Frame');
```

```
r.dynamics.mean = mirmean(r.dynamics.tmp.rms);
```

```
r.dynamics.lowenergy = mirlowenergy(r.dynamics.tmp.rms, 'ASR');
```

*r.timbral* = ***mirst****struct*;

...

*mireval*(*r*, 'Folders');

In this case, the temporary variable (here, *r.dynamics.tmp.rms*) should be used only inside that particular field in which it has been defined (here, *r.dynamics*).

# REFERENCES

- Alluri, V. & Toiviainen, P. (2010). Exploring perceptual and acoustic correlates of polyphonic timbre. *Music Perception*, 27(3), 223–241.
- Alluri, V., Toiviainen, P., Jääskeläinen, I., Sams, M., Glerean, E., & Brattico, E. (2012). Large-scale brain networks emerge from dynamic processing of musical timbre, key and rhythm. *Neuroimage* 59, 3677–3689.
- Alonso, David, Richard. A study of tempo tracking algorithms from polyphonic music signals, 4 COST 276 Workshop, 2003.
- Juan P. Bello, Chris Duxbury, Mike Davies, and Mark Sandler, On the Use of Phase and Energy for Musical Onset Detection in the Complex Domain, *IEEE Signal Processing Letters*, 11-6, 2004.
- Paul Boersma (1993). [Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound](#). *IFA Proceedings* 17: 97–110.
- Paul Boersma & David Weenink (2005). *Praat: doing phonetics by computer* [Computer program] <http://www.praat.org/>
- B. P. Bogert, M. J. R. Healy, and J. W. Tukey: The quefrency alalysis of time series for echoes: cepstrum, pseudo-autocovariance, cross-cepstrum, and saphe cracking. *Proceedings of the Symposium on Time Series Analysis* (M. Rosenblatt, Ed) Chapter 15, 209–243. New York: Wiley, 1963.
- T. Eerola, O. Lartillot, P. Toiviainen, "Prediction of Multidimensional Emotional Ratings in Music From Audio Using Multivariate Regression Models", [International Conference on Music Information Retrieval](#), Kobe, 2009.
- H. Fastl. Fluctuation strength and temporal masking patterns of amplitude-modulated broad-band noise. *Hearing Research*, 8:59–69, 1982.
- Y. Feng, Y. Zhuang, Y. Pan. (2003) Popular music retrieval by detecting mood, *ACM SIGIR Conference on Research and Development in Information Retrieval*.
- J. Foote, M. Cooper, U. Nam. (2002) "Audio Retrieval by Rhythmic Similarity", *ISMIR 2002*.
- Foote, Cooper. (2003). Media Segmentation using Self-Similarity Decomposition. *SPIE Storage and Retrieval for Multimedia Databases*, 5021, 167–75.
- Gómez, E. (2006). *Tonal description of music audio signal*. Phd thesis, Universitat Pompeu Fabra, Barcelona .
- Harte, C. A. and M. B. Sandler, Detecting harmonic change in musical audio, in *Proceedings of Audio and Music Computing for Multimedia Workshop*, Santa Barbara, CA, 2006.
- A. Klapuri, "Sound onset detection by applying psychoacoustic knowledge", *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 1999.
- Klapuri, A., A. Eronen and J. Astola. (2006). "Analysis of the meter of acoustic musical signals", *IEEE Transactions on Audio, Speech and Langage Processing*, 14-1, 342– 355.
- R. Koch, A Tutorial in Probability and Statistics, [http://web.cecs.pdx.edu/~roy/tutorial/Stat\\_int.htm](http://web.cecs.pdx.edu/~roy/tutorial/Stat_int.htm)
- Krimphoff, J., McAdams, S. & Winsberg, S. (1994), Caractérisation du timbre des sons complexes. II : Analyses acoustiques et quantification psychophysique. *Journal de Physique*, 4(C5), 625–628.
- Krumhansl, *Cognitive foundations of musical pitch*. Oxford UP, 1990.
- K. Jacobson, *A multifaceted approach to music similarity*, ISMIR 2006.
- Jensen, *Timbre Models of Musical Sounds*, Rapport 99/7, University of Copenhagen, 1999.
- Juslin, P. N. (2000). Cue utilization in communication of emotion in music performance: relating performance to perception. *Journal of Experimental Psychology: Human Perception and Performance*, 26(6), 1797–813.
- Lartillot, O., Eerola, T., Toiviainen, P., Fornari, J., "Multi-feature modeling of pulse clarity: Design, validation, and optimization", [International Conference on Music Information Retrieval](#), Philadelphia, 2008.



- Lartillot, O., Cereghetti, D., Eliard, K., Grandjean, D., "A simple, high-yield method for assessing structural novelty", [3rd International Conference on Music & Emotion](#), Jyväskylä, 2013.
- Lartillot, Olivier & Grandjean, Didier (2019). Tempo and Metrical Analysis by Tracking Multiple Metrical Levels Using Auto-correlation. *Applied Sciences*, 9(23). doi: 10.3390/app9235121.
- Olivier Lartillot, Kristian Nymoen, Guilherme Schmidt Câmara, & Anne Danielsen (2021). Computational localization of attack regions through a direct observation of the audio waveform, *Journal of the Acoustical Society of America (JASA)*, 149(1), 723—736.
- Lartillot, O. & Weisser, S. (2021). Roughness, Crackliness, Buzzingness, ...: Characterizations of Sonic Unsteadiness and Application to the Analysis of Traditional Music from Ethiopia, Kenya, Morocco and India. *16th International Conference on Music Perception and Cognition (ICMPC) & the 11th triennial conference of ESCOM*.
- Laukka, P., Juslin, P. N., and Bresin, R. (2005). A dimensional approach to vocal expression of emotion. *Cognition and Emotion*, 19, 633-653.
- Logan, B., A. Salomon. (2001). A content-based music similarity function. Cambridge Research Laboratory, Technical Report Series.
- McAulay, R.; Quatieri, T. (1996). "Speech analysis/Synthesis based on a sinusoidal representation", *Acoustics, Speech and Signal Processing, IEEE Transactions on*, Volume 34, Issue 4. Page(s): 744 - 754
- Nabney, I. "NETLAB: Algorithms for pattern recognition", *Springer Advances In Pattern Recognition Series*, 2002.
- Nymoen, K., A. Danielsen and J. London. "Evaluating Matlab toolboxes for estimating attack phase descriptors from audio files", *14th Sound and Music Computing Conference*, 2017.
- Pampalk, E. "A Matlab Toolbox to Compute Similarity from Audio", *International Conference on Music Information Retrieval*, 2004
- E. Pampalk, A. Rauber, D. Merkl, "Content-based Organization and Visualization of Music Archives", *ACM Multimedia* 2002, pp. 570-579.
- R. D. Patterson et al. "Complex sounds and auditory images," in *Auditory Physiology and Perception*, Y. Cazals et al, Oxford, 1992, pp. 429-446.
- Peeters. G. (2004). *A large set of audio features for sound description (similarity and classification) in the CUIDADO project*. version 1.0
- Peeters, G. Music pitch representation by periodicity measures based on combined temporal and spectral representations. *ICASSP* 2006.
- Peeters, G., B.L. Giordano, P. Susini, N. Misdariis, S. McAdams. "The Timbre Toolbox: Extracting audio descriptors from musical signals", *Journal of the Acoustical Society of America*, 2902-16, 2011.
- Plomp & Levelt "Tonal Consonance and Critical Bandwidth" *Journal of the Acoustical Society of America*, 1965.
- Pohle, T., E. Pampalk and G. Widmer (2005). Evaluation of Frequently Used Audio Features for Classification of Music Into Perceptual Categories. *Proceedings of the Fourth International Workshop on Content-Based Multimedia Indexing (CBMI'05)*, Riga, Latvia, June 21-23.
- Scheirer, E. D. (1998). "Tempo and beat analysis of acoustic musical signals", *Journal of the Acoustical Society of America*, 103-1, 588-601.
- Sethares, W. A. (1998). *Tuning, Timbre, Spectrum, Scale*, Springer-Verlag.
- Shannon, C. E. 1948. "A Mathematical Theory of Communication." *Bell Systems Technical Journal* 27:379-423.
- Slaney, M. *Auditory Toolbox Version 2*, Technical Report. Interval Research Corporation, 1998-010, 1998.
- S.S. Stevens and J. Volkman (1940) "The relation of pitch to frequency: A revised scale" *Am. J. Psychol.* 53: 329-353.
- Terhardt, E. Calculating virtual pitch. *Hearing Research*, 1:155-182, 1979.
- Toiviainen & Krumhansl, "Measuring and modeling real-time responses to music: The dynamics of tonality induction", *Perception* 32-6, pp. 741-766, 2003.

- Toiviainen, Snyder. "Tapping to Bach: Resonance-based modeling of pulse", *Music Perception*, 21-1, 43-80, 2003.
- Tolonen, Karjalainen. A Computationally Efficient Multipitch Analysis Model, *IEEE Transactions on Speech and Audio Processing*, 8(6), 2000.
- Tzanetakis, Cook. Musical genre classification of audio signals. *IEEE Tr. Speech and Audio Processing*, 10(5), 293-302, 2002.
- Van Noorden, L., & Moelants, D. (1999). Resonance in the Perception of Musical Pulse. *Journal of New Music Research*, 28(1), 43-66.
- Vassilakis, P. N. (2001). Perceptual and Physical Properties of Amplitude Fluctuation and their Musical Significance. Doctoral Dissertation. Los Angeles: University of California, Los Angeles; Systematic Musicology.
- Vesanto, J. "Self-Organizing Map in Matlab: the SOM Toolbox", *Matlab DSP Conference*, 35-40, 1999.
- Weisser, S. & Lartillot, O. "Investigating non-Western musical timbre: a need for joint approaches", *3rd International Workshop on Folk Music Analysis*, Amsterdam, 2013.
- E. Zwicker, G. Flottorp and S.S. Stevens (1957) "Critical bandwidth in loudness summation" *J. Acoust. Soc. Am.* 29: 548-557.