

Audio Content Description

Mathieu Lagrange 



October 17, 2018

Outline

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



Outline

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



Outline

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



Outline

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



Outline

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



Outline

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks

Perceptual vs. physical properties

- ① Loudness \Leftrightarrow Amplitude
- ② Pitch \Leftrightarrow Fundamental frequency
- ③ Timbre \Leftrightarrow spectrum and envelope ??

Perceptual vs. physical properties

- ① Loudness \Leftrightarrow Amplitude

Definition

Loudness is that attribute of auditory sensation in terms of which sounds can be ordered on a scale extending from quiet to loud

- ② Pitch \Leftrightarrow Fundamental frequency
- ③ Timbre \Leftrightarrow spectrum and envelope ??

Perceptual vs. physical properties

- ① Loudness \Leftrightarrow Amplitude
- ② Pitch \Leftrightarrow Fundamental frequency

Definition

Pitch is a perceptual property that allows the ordering of sounds on a frequency-related scale

- ③ Timbre \Leftrightarrow spectrum and envelope ??

Perceptual vs. physical properties

- ① Loudness \Leftrightarrow Amplitude
- ② Pitch \Leftrightarrow Fundamental frequency
- ③ Timbre \Leftrightarrow spectrum and envelope ??

Definition

Timbre is what makes a particular musical sound different from another, even when they have the same pitch and loudness

Musical properties

Music is about organization of sounds in time

- ⌚ Dynamics: *pp* to *ff*
- ⌚ Duration: ♩ ♫ .
- ⌚ Rhythm: ♪

Definition

rhythm is the timing of events on a human scale

- ⌚ Beat

Definition

Beat is the rhythm listeners would tap their toes to when listening to a piece of music

Timbre

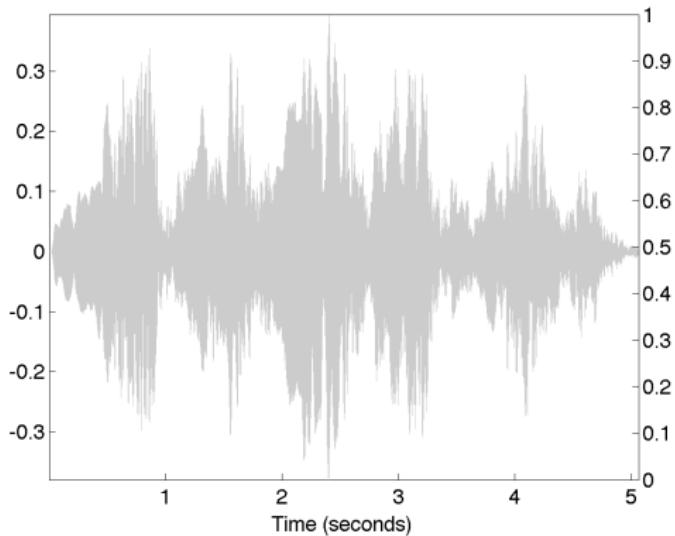
Five major acoustic parameters

- ① The range between tonal and noiselike character
- ② The spectral envelope
- ③ The time envelope in terms of rise, duration, and decay
- ④ The changes both of spectral envelope (formant-glide) and fundamental frequency (micro-intonation)
- ⑤ The prefix, or onset of a sound, quite dissimilar to the ensuing lasting vibration

Female singing voice



↪ waveform



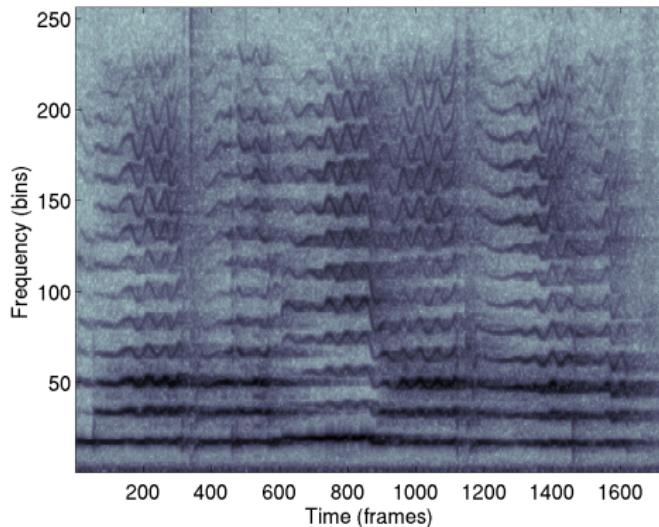
↪ spectrogram

↪ both

Female singing voice



- waveform
- spectrogram

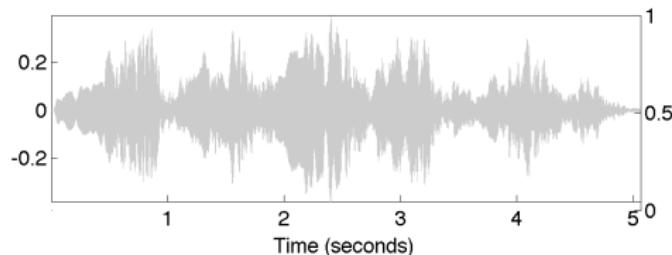
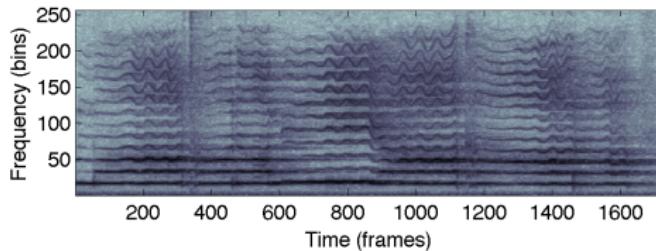


- both

Female singing voice



- waveform
- spectrogram
- both



① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks

Loudness

Loudness is the perceptual attribute of sound that correlate with its physical strength

ε a useful scale is the deciBels:

$$v_{dB} = 20 \log_{10} \left(\frac{v(n)}{v_0} \right)$$

ε v_0 is a reference constant, set to 1 when $-1 < v_n < 1$

Root Mean Square (RMS)

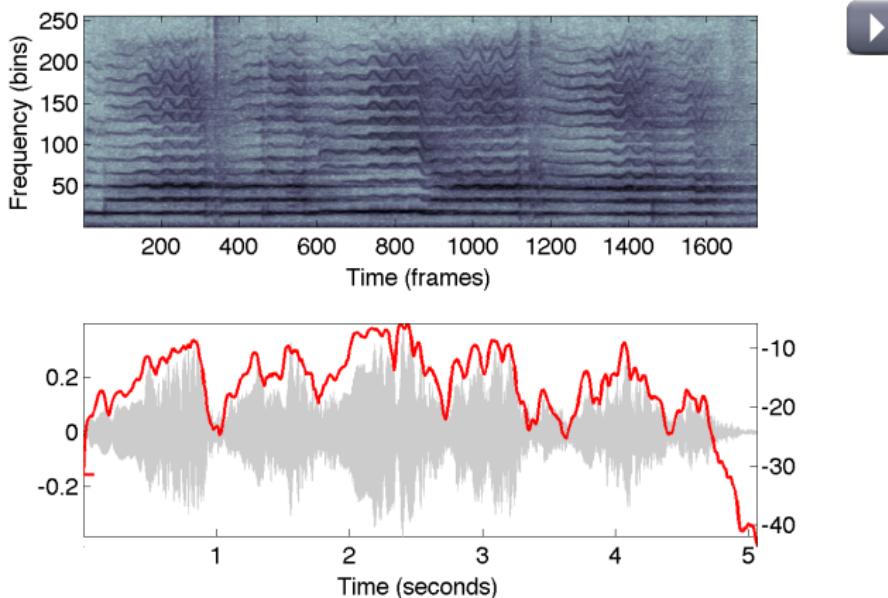


Definition

The Root Mean Square (RMS) is defined as

$$v(n) = \sqrt{\frac{1}{N} \sum x(i)^2}$$

Root Mean Square (RMS)



- ⌚ computation per blocks of hundred of milliseconds
- ⌚ low values ⇔ low volume
- ⌚ large values ⇔ high volume

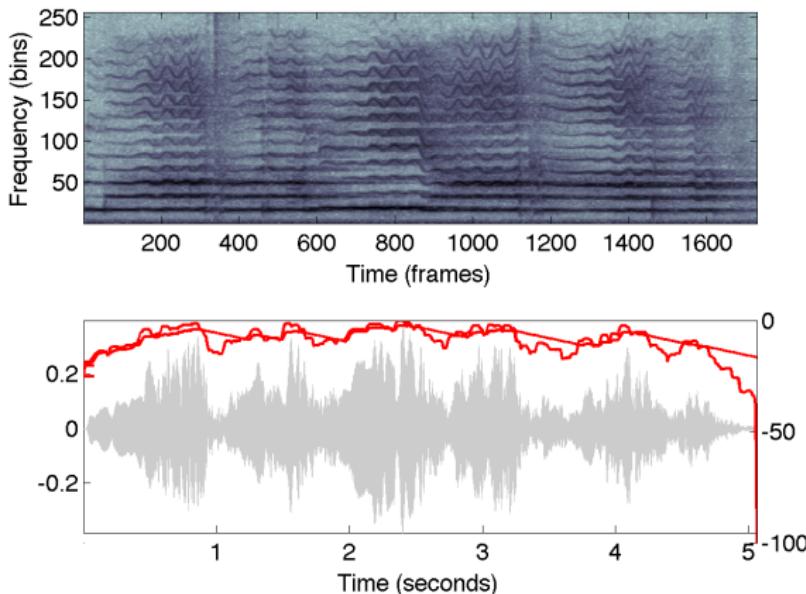
Peak Envelope



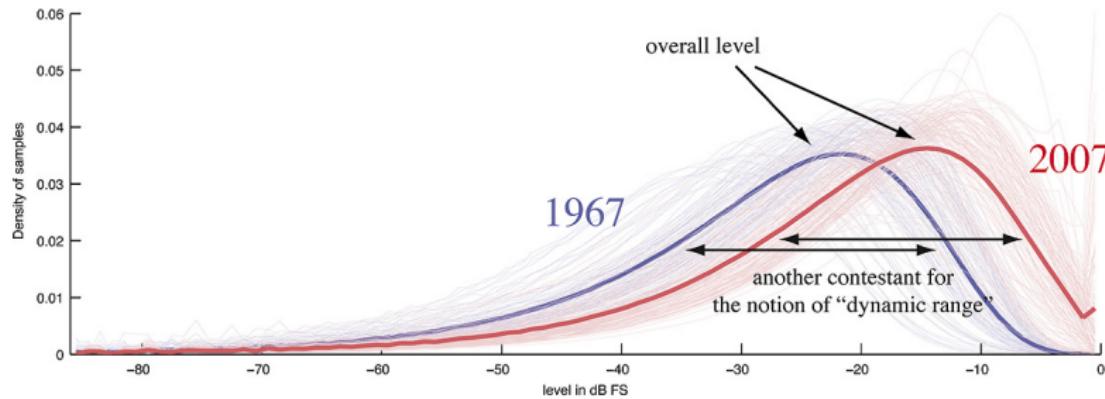
Definition

The Peak Envelope is defined as the maximal value within a given block of samples.

Peak Envelope



The loudness war



1967



2007



Comparisons between signal levels and picture levels as defined in Photoshop result in another interpretation of the loudness war.

Emmanuel Deruty Loudness study of 4.5K tracks

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks

Pre-processing

- ε Down mixing: average all channels contribution
- ε DC removal: subtract arithmetic mean
- ε Normalization: set (maximal / averaged) amplitude to a predefined value

Pre-processing

- ε Down mixing: average all channels contribution

$$x(i) = \frac{1}{C} \sum_{c=1}^C x_c(i)$$

- ε DC removal: subtract arithmetic mean
- ε Normalization: set (maximal / averaged) amplitude to a predefined value

Pre-processing

- ⌚ Down mixing: average all channels contribution
- ⌚ DC removal: subtract arithmetic mean

$$x(i) = x(i) - \frac{1}{N} \sum_{n=1}^N x(n)$$

- ⌚ Normalization: set (maximal / averaged) amplitude to a predefined value

Pre-processing

- ⌚ Down mixing: average all channels contribution
- ⌚ DC removal: subtract arithmetic mean
- ⌚ Normalization: set (**maximal** / averaged) amplitude to a predefined value

$$x(i) = \frac{x(i)}{\max(|x(i)|)}$$



Pre-processing

- ⌚ Down mixing: average all channels contribution
- ⌚ DC removal: subtract arithmetic mean
- ⌚ Normalization: set (maximal / averaged) amplitude to a predefined value

$$x(i) = \frac{x(i)}{\sqrt{\sum_{n=1}^N x(n)^2}}$$



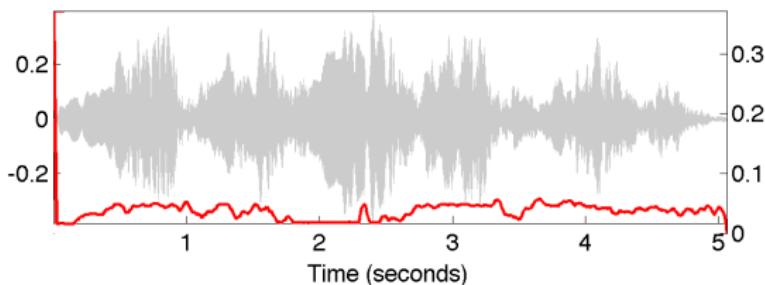
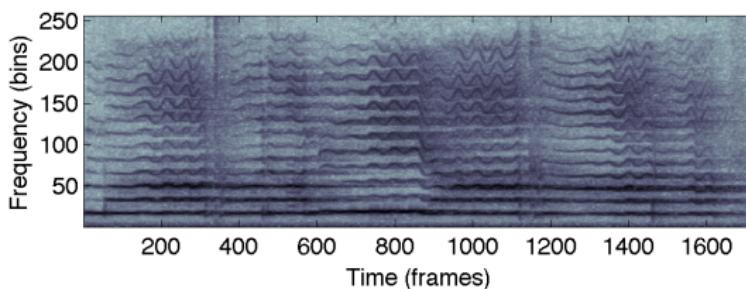
The zero crossing rate (ZCR)



Definition

The number of changes of sign of audio samples in consecutive block.

The zero crossing rate (ZCR)



- ⌚ low values \Leftrightarrow low frequency content
- ⌚ large values \Leftrightarrow high frequency content
- ⌚ large value change across blocks indicates noisy content

The predictivity ratio

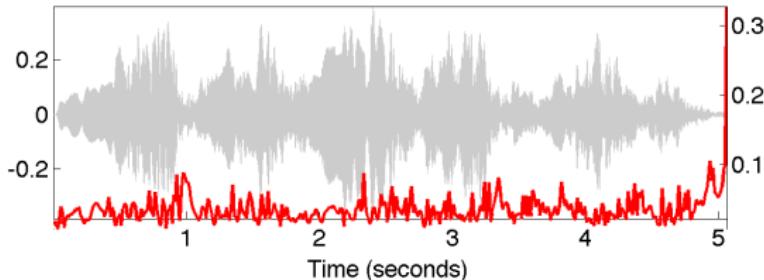
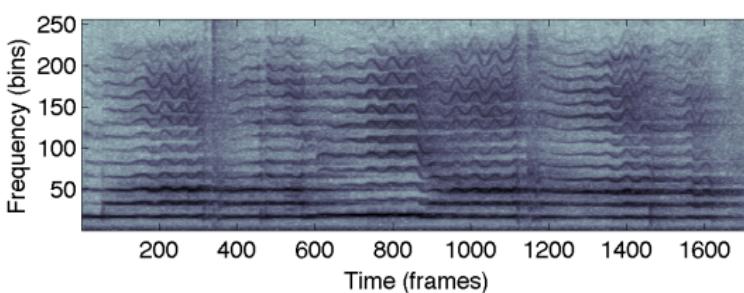


Definition

The predictivity ratio is a measure of how well the audio signal can be predicted by an K-order linear prediction:

$$\hat{x}(i) = \sum_{k=1}^K b_k x(i-k)$$

The predictivity ratio



- ⌚ low values \Leftrightarrow tonal signal
- ⌚ large values \Leftrightarrow unpredictable signals (noise)

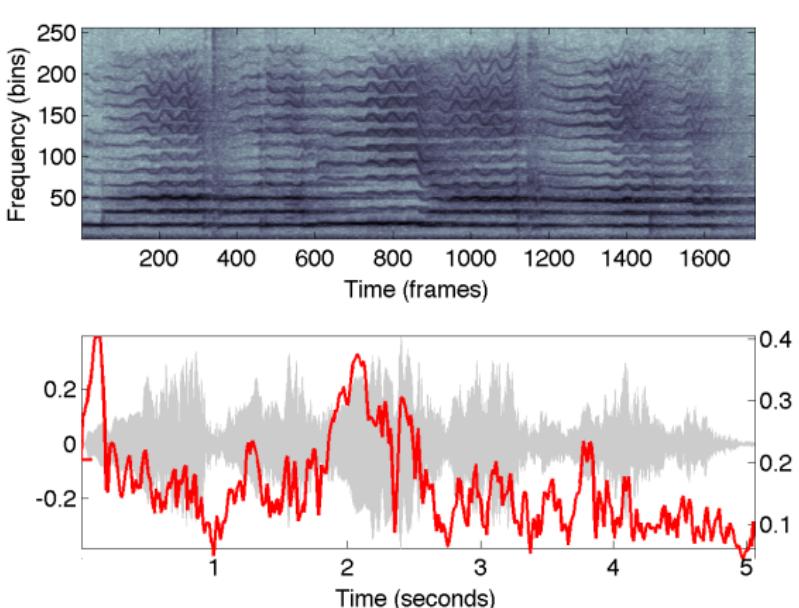
The spectral crest factor



Definition

The spectral crest factor is the ratio of the maximum in the magnitude spectrum with the sum of the magnitude spectral bins.

The spectral crest factor



- ⌚ low values \Leftrightarrow flat magnitude spectrum
- ⌚ large values \Leftrightarrow sinusoids

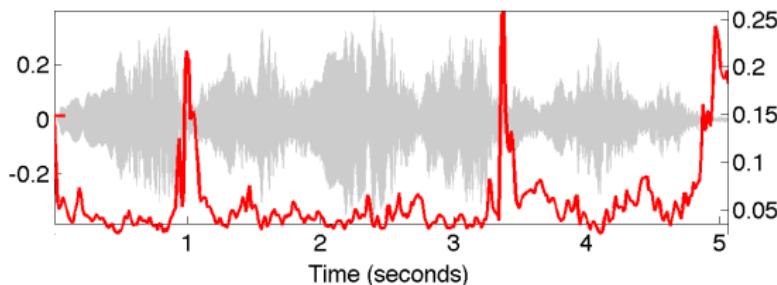
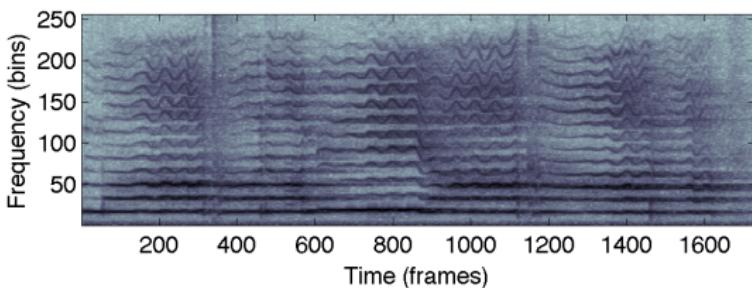
The spectral flatness



Definition

The spectral flatness is the ratio of geometric mean and arithmetic mean of the magnitude spectrum.

The spectral flatness



- ⌚ low values \Leftrightarrow possibly tonal components
- ⌚ large values \Leftrightarrow flat (noisy) spectrum

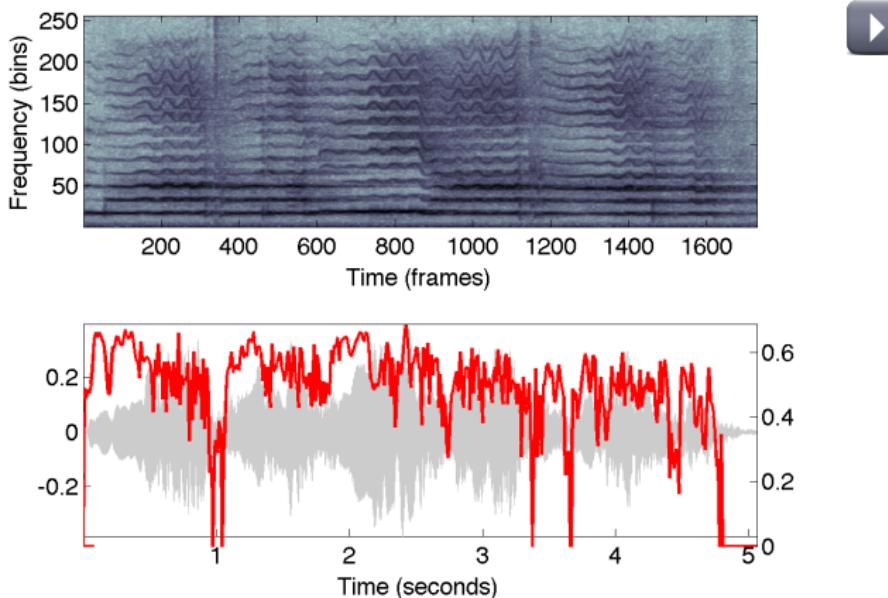
The tonal power ratio



Definition

The tonal power ratio is the ratio between the tonal power and overall power.

The tonal power ratio



Tonal components are

- ε local maxima: $|X(k-1, n)|^2 < |X(k, n)|^2 > |X(k+1, n)|^2$
- ε powerful: magnitude above a given threshold

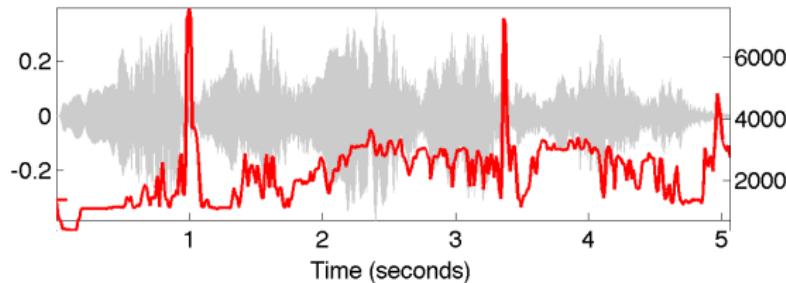
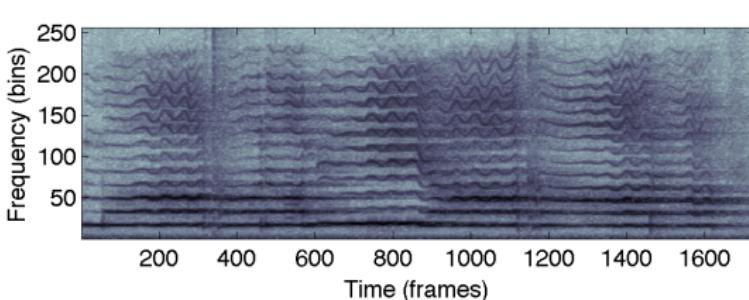
The spectral rolloff



Definition

The frequency bin below which the accumulated magnitude of the spectrum reach a certain percentage of the overall magnitude spectrum.

The spectral rolloff



- ⌚ low values \Leftrightarrow low audio bandwidth
- ⌚ large values \Leftrightarrow noise

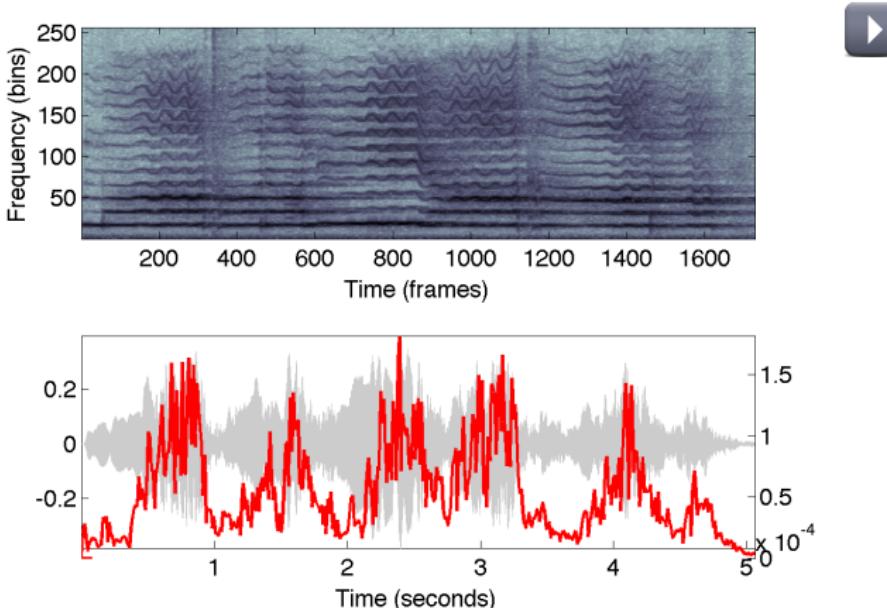
The spectral flux (roughness)



Definition

The spectral flux is the average difference between consecutive magnitude bins.

The spectral flux (roughness)



- ⌚ low values \Leftrightarrow steady state signals or low input levels
- ⌚ large values \Leftrightarrow abrupt changes

The spectral centroid (brightness / sharpness)

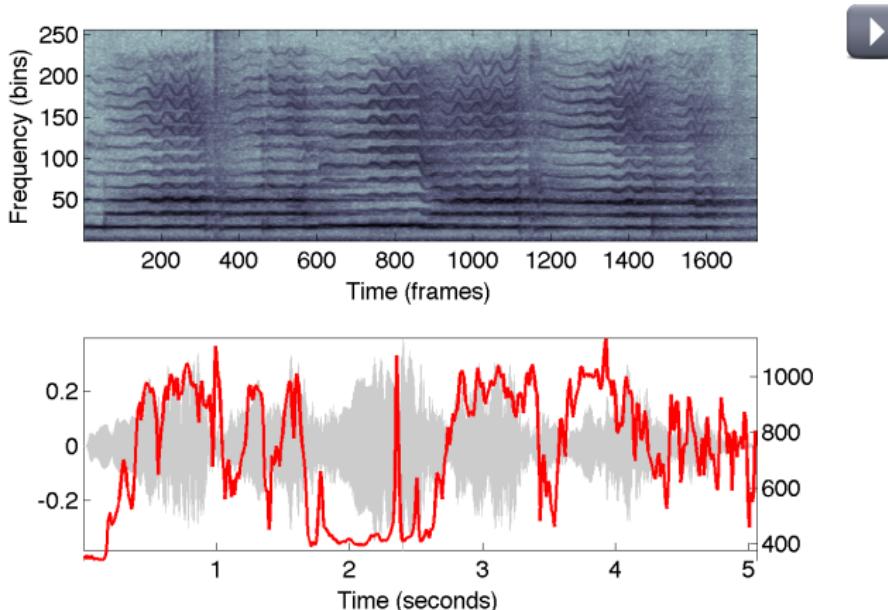


Definition

The spectral centroid is defined as the Center of Gravity (COG) of the spectral energy:

$$sc(n) = \frac{\sum_{k=1}^K k|X(k,n)|^2}{\sum_{k=1}^K |X(k,n)|^2}$$

The spectral centroid (brightness / sharpness)



- ⌚ low values \Leftrightarrow dominant low frequencies
- ⌚ large values \Leftrightarrow dominant high frequencies

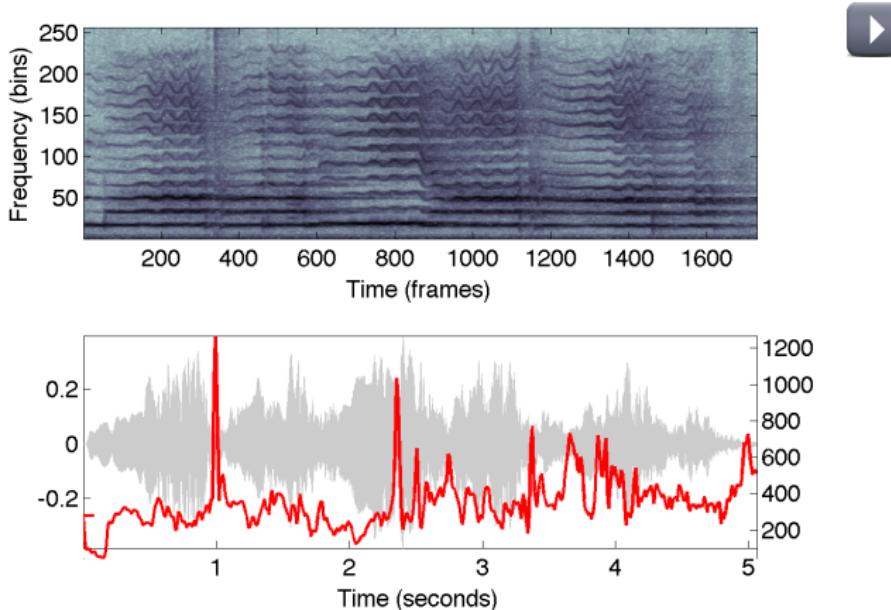
The spectral spread (instantaneous bandwidth)



Definition

The spectral spread describes the concentration of the spectral energy around the spectral centroid.

The spectral spread (instantaneous bandwidth)



- ⌚ low values \Leftrightarrow during notes (monophonic case)
- ⌚ large values \Leftrightarrow during pauses and transients

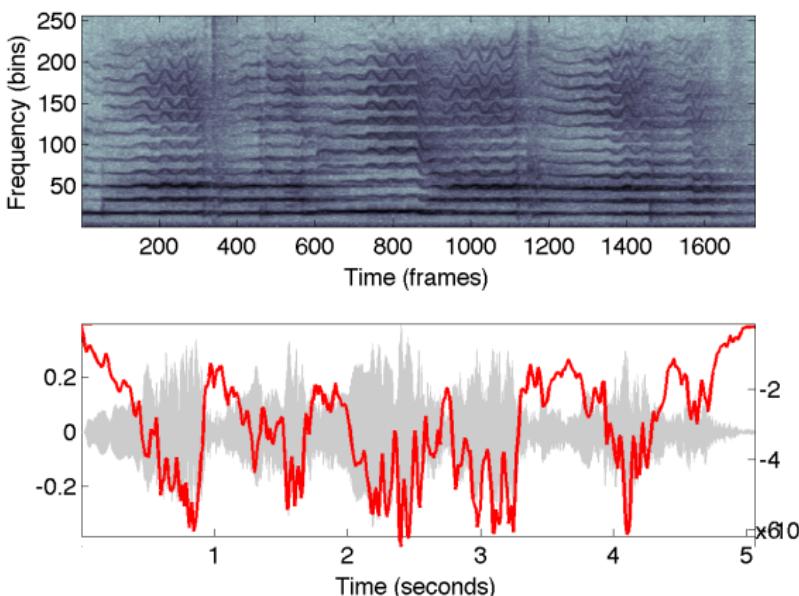
The spectral slope



Definition

The spectral slope is computed from a linear regression over frequency of the magnitude spectrum.

The spectral slope



- ⌚ low values \Leftrightarrow "flat" spectrum
- ⌚ large values \Leftrightarrow few higher harmonics

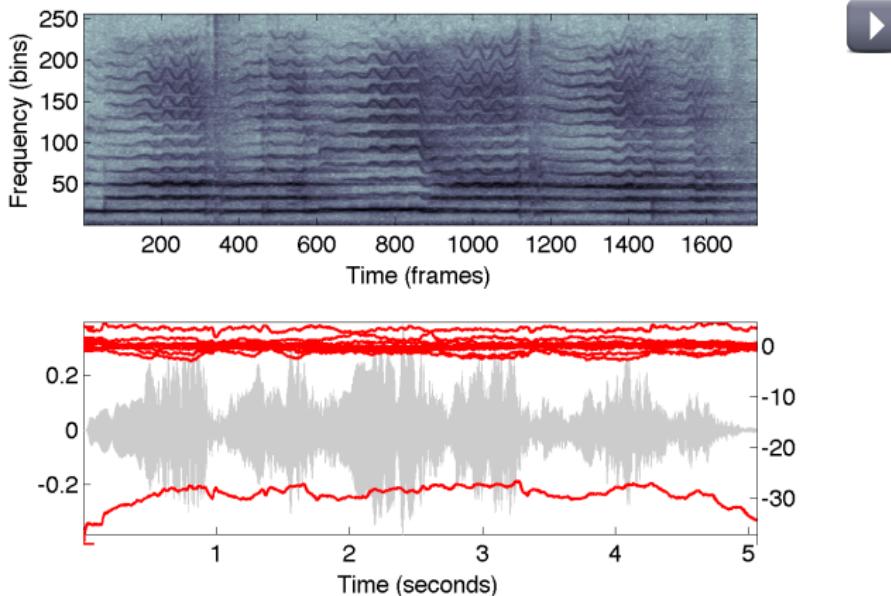
Mel Frequency Cepstral Coefficients (MFCC)s



Definition

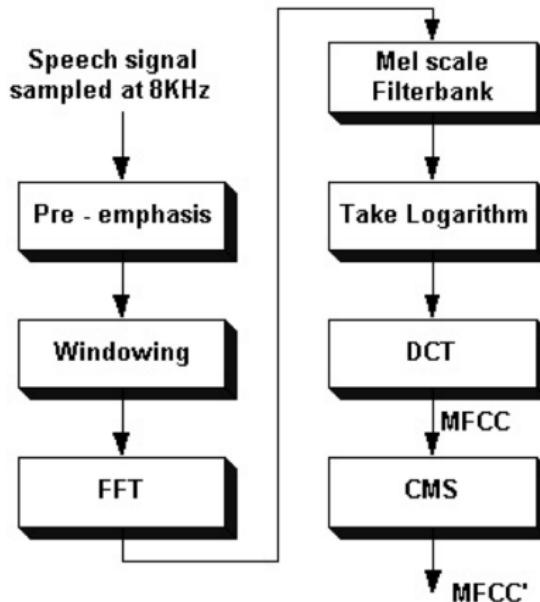
The Mel Frequency Cepstral Coefficients (MFCC)s can be seen as a compact description of the shape of the spectral envelope of an audio signal.

Mel Frequency Cepstral Coefficients (MFCC)s



- ⌚ based on a mel-warped spectrum
- ⌚ a logarithm of the magnitude spectrum
- ⌚ apply a DCT on the resulting spectral representation

Mfccs



Mel Frequency Cepstral Coefficients (MFCC)s

- ① Take the Fourier transform of (a windowed excerpt of) a signal.
- ② Map the powers of the spectrum obtained above onto the mel scale, using overlapping windows.
- ③ Take the logs of the powers at each of the mel frequencies.
- ④ Take the discrete cosine transform (DCT)
- ⑤ The MFCCs are the amplitudes of the resulting spectrum.

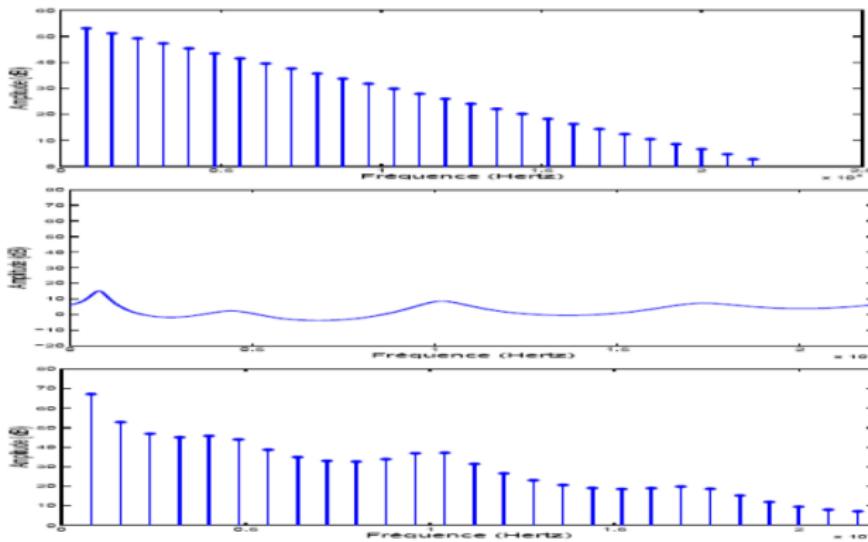
Mfccs : log spectrum



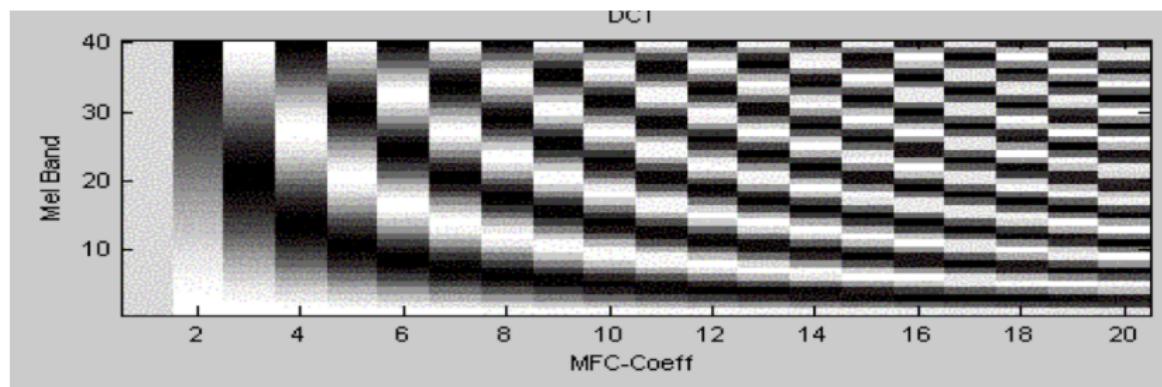
linear frequency

log frequency

Mfccs : source / filter model



Mfccs : dct basis



Projecting the magnitude spectrum over this basis allows us to capture the envelope with the first resulting dct coefficients.

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks

Pitch

Many tonal sounds produced by many musical instruments

- ε can be considered as a sum of **sinusoidal** components with frequencies that are **harmonically** related: $f_0, 2f_0, 3f_0, \dots$
- ε then the **fundamental frequency f_0** dominates the pitch perception
- ε mostly a relative perception as the fundamental frequency may be missing

Scales

The relation between the fundamental frequency and the perceived pitch is non-linear

- ε this non-linearity is tied to the frequency resolution of the human cochlea
- ε Mel scale: the scale is a measure of tone height. One analytic model:

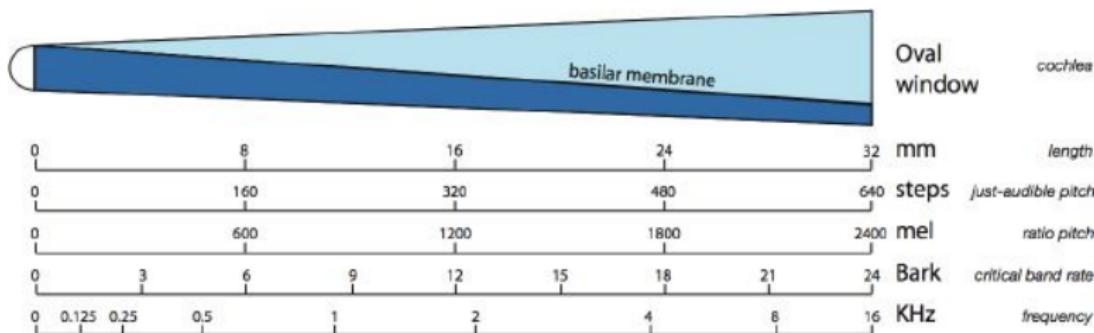
$$m_F(f) = 1000 \log_2 \left(1 + \frac{f}{1000} \right)$$

- ε Bark scale: the scale come from the bandwidth of measured critical bands:

$$b_F(f) = 7 \operatorname{arcsinh} \left(\frac{f}{650} \right)$$

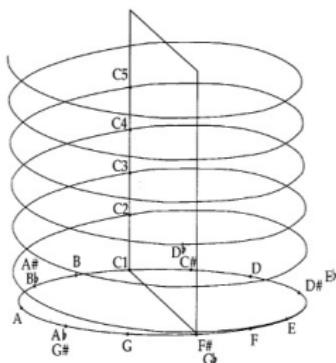
- ε the 2 scales are related by 1 bark = 100 mel

Scales



Physical relationships between the different scales

Chroma perception



Humans also tends to group pitches with ratios of 2.

The frequency of musical pitch

The MIDI standard

- ε uses the equal temperament which results in enharmonic equivalence
- ε each semi-tone has a distance of 1 to its nearest neighbor
- ε from MIDI pitch to frequency:

$$p(f) = 69 + 12 \log_2 \left(\frac{f}{f_{A4}} \right)$$

- ε from frequency to MIDI pitch:

$$f(p) = f_{A4} 2^{\frac{p-69}{12}}$$

Tuning frequency

The tuning frequency f_{A4}

- ⌚ is the frequency of the **concert** pitch or **standard** pitch
- ⌚ is standardized internationally to 440 Hz
- ⌚ that said the exact frequency used by musicians can vary due to the use historic instruments or timbre preferences

Tuning frequency

The tuning frequency f_{A4}

- ⌚ is the frequency of the **concert** pitch or **standard** pitch
- ⌚ is standardized internationally to 440 Hz
- ⌚ that said the exact frequency used by musicians can vary due to the use historic instruments or timbre preferences

Tuning frequency

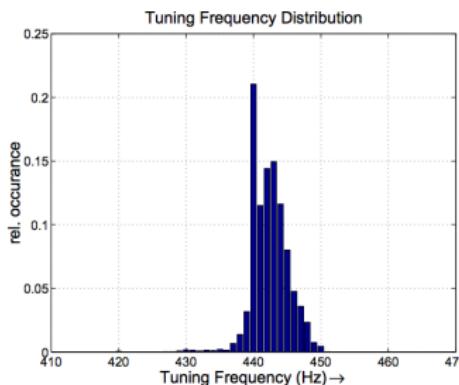
The tuning frequency f_{A4}

- ⌚ is the frequency of the **concert** pitch or **standard** pitch
- ⌚ is standardized internationally to 440 Hz
- ⌚ that said the exact frequency used by musicians can vary due to the use historic instruments or timbre preferences

Tuning frequency

The tuning frequency f_{A4}

- ⌚ is the frequency of the **concert** pitch or **standard** pitch
- ⌚ is standardized internationally to 440 Hz
- ⌚ that said the exact frequency used by musicians can vary due to the use historic instruments or timbre preferences



Identification of tonal components

For a wide range of applications (coding, analysis, synthesis), identifying the tonal components is of crucial interest.

The input signal can be assumed to be:

- ⌚ a time variant mixture of tonal and non-tonal components
- ⌚ an undefined number of tonal voices

The tonal components can be assumed to be

- ⌚ salient: with a certain intensity and not masked by nearby components
- ⌚ deterministic: its phase cannot change erratically
- ⌚ stationary for a minimum length of time



Methods

- ⌚ Local maximum: binary decision
- ⌚ Peakiness: magnitude difference between the local maxima and neighbors
- ⌚ Thresholding: adaptive relative to a smoothed magnitude spectrum
- ⌚ Frequency coherence: analysis of the phase

Fundamental frequency detection

Time domain methods:

- ⌚ Zero Crossing Rate (ZCR): the interval between 2 zero crossings related directly to the fundamental period length
- ⌚ Autocorrelation function (ACF): the goal of the autocorrelation is to correlate the signal and a shifted version of itself. The mathematical definitions is:

$$y(n) = \sum_{k=1}^M x(n)x(n+k)$$

If the signal is periodic, the autocorrelation function $y(n)$ also will be, and if the signal is harmonic the autocorrelation function will have peaks in multiples of the fundamental frequency. This technique is most efficient at mid to low frequencies.



Fundamental frequency detection

Frequency domain methods:

- ε Harmonic Product Spectrum:

$$X(k, n) = \prod_{j=1}^O |X(jk, n)|^2$$

- ε Harmonic Sum Spectrum:

$$X(k, n) = \sum_{j=1}^O |X(jk, n)|^2$$

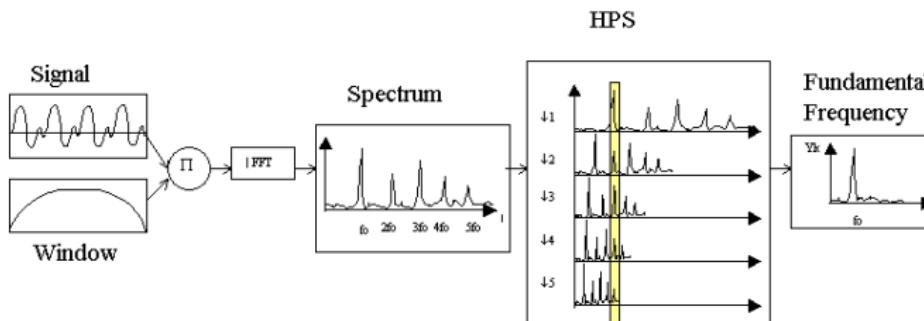
- ε ACF of the magnitude spectrum
- ε Cepstral domain

Fundamental frequency detection

Frequency domain methods:

- ↪ Harmonic Product Spectrum:

$$X(k, n) = \prod_{j=1}^O |X(jk, n)|^2$$



- ↪ Harmonic Sum Spectrum:

$$X(k, n) = \sum_{j=1}^O |X(jk, n)|^2$$

Pitch Chroma

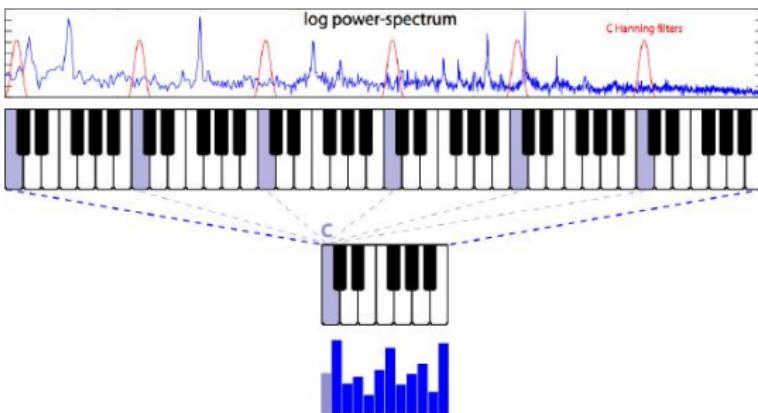
The pitch chroma

is the histogram-like 12-dimensional vector with each dimension representing one pitch class.

The main building blocks are:

- ⌚ frequency representation grouped into semi tone bands
- ⌚ a measure of salience in each band
- ⌚ computation of the sum of bands over all octaves

Pitch Chroma



This feature is largely invariant to:

- ⌚ octave change
- ⌚ timbre fluctuation
- ⌚ noise

① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks



Human perception of temporal events

- ε Onset: the start of a sound event defined as its time and strength
- ε Attack time: 5 to 200 ms depending on the musical instrument
- ε Tempo: rate at which different pulses with equal duration units occur at a moderate and natural rate. Perceived tempo is the **tactus** or the foot tapping rate (typical value is 100 Beats per Minute (BPM))

$$T = \frac{60}{t(j+1) - t(j)}$$

- ε Meter: alternation of strong and weak musical elements which are grouped with a length of normally 3 to 7 beats or of duration of around 5 seconds.
- ε Rhythm: perceptual groups of musical elements of length between one beat and the length of the meter

Human perception of temporal events

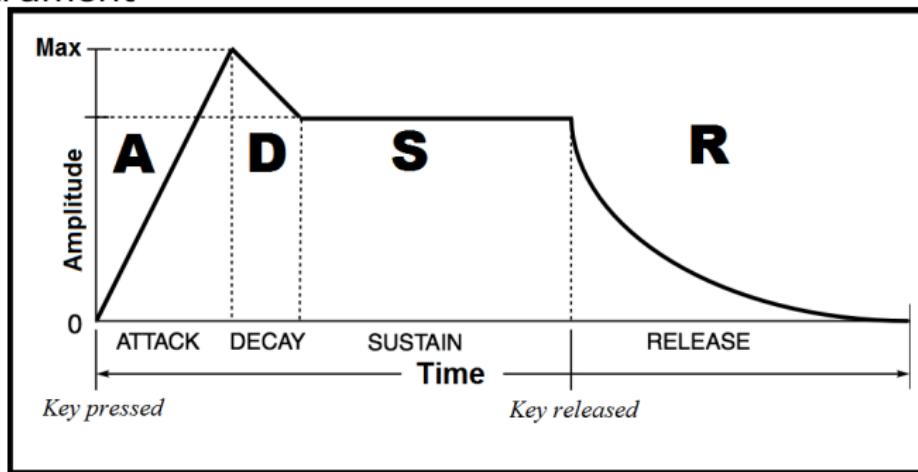
- ε Onset: the start of a sound event defined as its time and strength
- ε Attack time: 5 to 200 ms depending on the musical instrument
- ε Tempo: rate at which different pulses with equal duration units occur at a moderate and natural rate. Perceived tempo is the **tactus** or the foot tapping rate (typical value is 100 Beats per Minute (BPM))

$$T = \frac{60}{t(j+1) - t(j)}$$

- ε Meter: alternation of strong and weak musical elements which are grouped with a length of normally 3 to 7 beats or of duration of around 5 seconds.
- ε Rhythm: perceptual groups of musical elements of length between one beat and the length of the meter

Human perception of temporal events

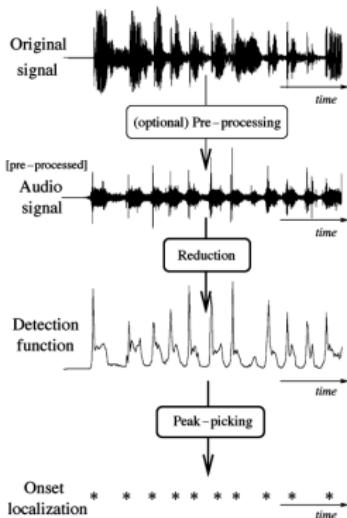
- Onset: the start of a sound event defined as its time and strength
- Attack time: 5 to 200 ms depending on the musical instrument



- Tempo: rate at which different pulses with equal duration units occur at a moderate and natural rate. Perceived tempo is the **tactus** or the foot tapping rate (typical value

Onset detection

Onset detection systems are usually built on 3 successive processing steps



Bello & al [A Tutorial on Onset Detection in Music Signals](#) IEEE transactions on Speech and Audio Processing, vol. 13, no. 5, September 2005

Novelty function

Something new is happening:

ε spectral flux:

$$n_f = \sum_{k=k_{min}}^{k_{max}} \sqrt{|X(n, l)|} - \sqrt{|X(n - 1, k)|}$$

ε logarithmic distance:

$$n_f = \sum_{k=k_{min}}^{k_{max}} \log_2 \frac{|X(n, l)|}{|X(n - 1, k)|}$$

ε complex difference:

$$n_f = \sum_{k=k_{min}}^{k_{max}} |X(n, l) - X(n - 1, k)|$$



Peak picking

ε Fixed threshold:

$$T = \lambda$$

ε Moving Average (MA) smoothing:

$$T(n) = \lambda + \sum b(j)n_f(n-j)$$

ε Median filtering

Beat histogram

The Beat histogram

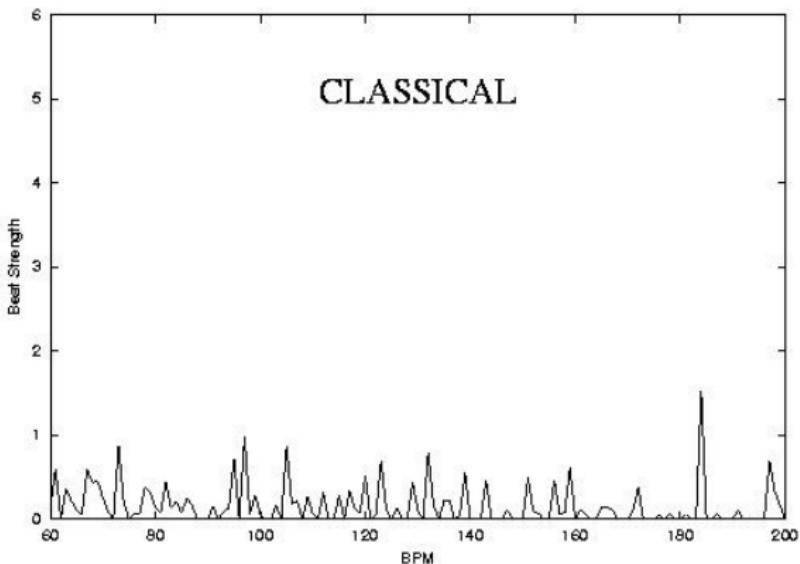
can be interpreted as the frequency domain representation of the novelty function. It is a way to visualize some rhythmic properties of the signal.

Typical processing blocks are:

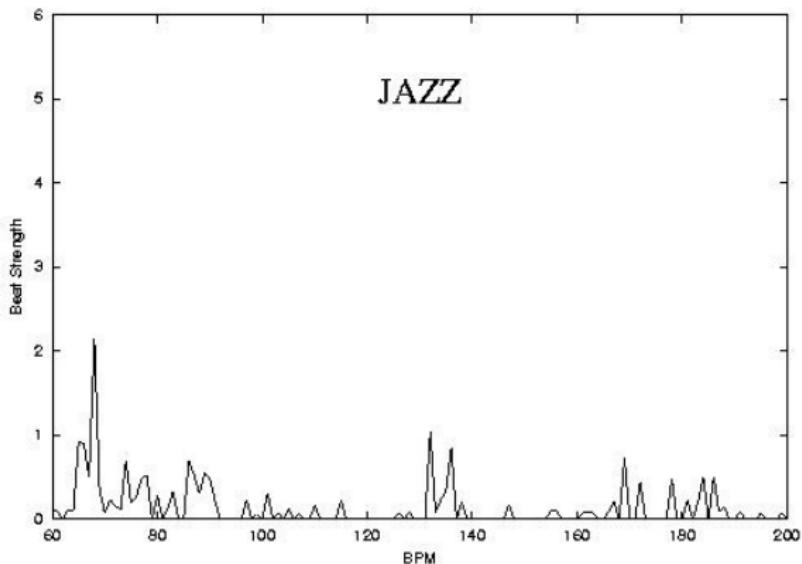
- ① split the signal into several frequency bands (4)
- ② Full Wave Rectification (FWR): compute absolute value
- ③ envelope smoothing: low pass filtering
- ④ down sampling: reducing the sampling rate
- ⑤ DC removal: subtracting the arithmetic mean
- ⑥ frequency analysis: ACF
- ⑦ amplitude thresholding: select the 3 most prominent components and add them to the histogram



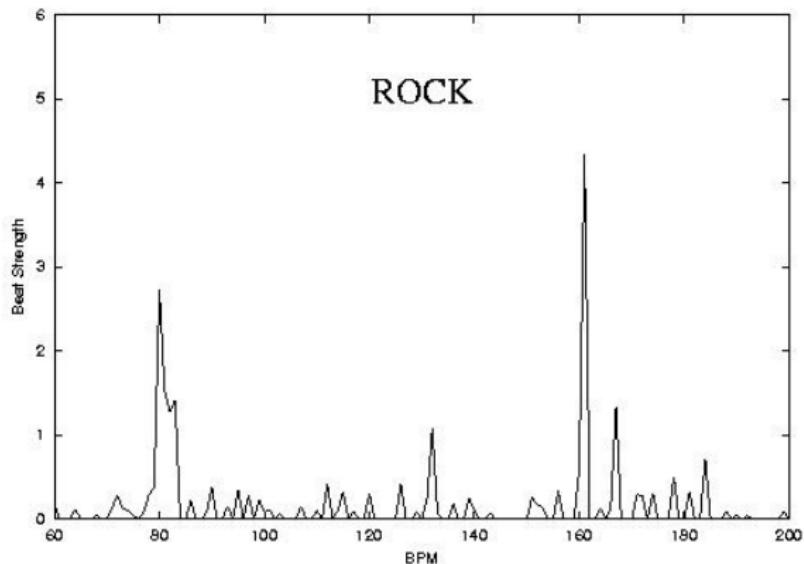
Beat histogram



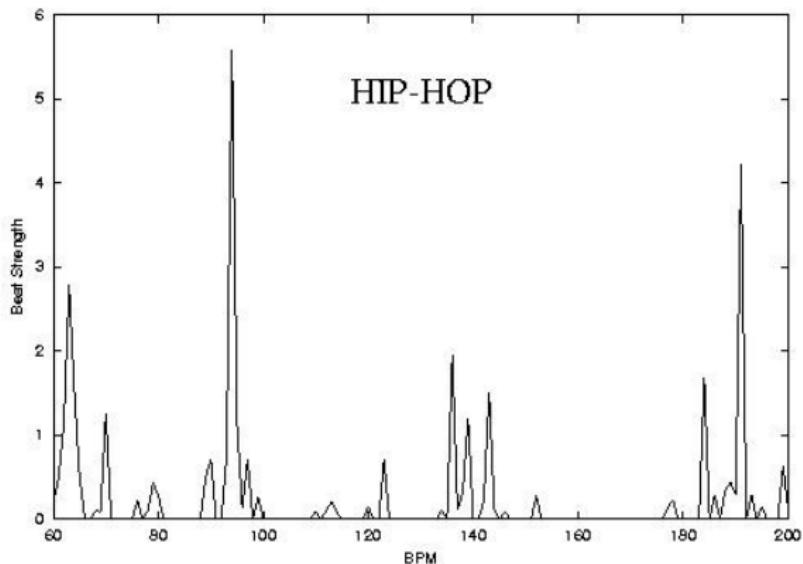
Beat histogram



Beat histogram



Beat histogram



① Properties

② Intensity

③ Spectral features

④ Tonal Analysis

⑤ Temporal Analysis

⑥ Tasks

Common machine listening tasks

Machine listening tasks

- ε Visualization
- ε Query by example
- ε Classification

For each of those tasks, the features are usually not considered directly.

Derived features

⌚ Time derivative:

$$v_{j,\Delta} = v_j(n) - v_j(n-1)$$

⌚ Low pass filtered version: to focus on long term variations

⌚ Monomials:

$$v_{jl} = v_j(n) \cdot v_l(n)$$



Normalization

- ↪ if the features can reasonably be considered as symmetric and similar shape (Gaussian distributed), the most standard way is the standardization

$$v_{j,N} = \frac{v_j(n) - \mu_{v_j}}{\sigma_{v_j}}$$

- ↪ if not, one can use some features transform such as the Box-Cox transform:

$$v_\lambda = \frac{v^\lambda - 1}{\lambda}$$

- ↪ or consider the median



Pooling

For classification tasks, it is common to smooth out local variations by considering a **texture window** of several frames. The pooling operator is:

- ε image processing: max
- ε audio processing: sum
- ε typical window length ranges from 1 to 15 seconds
- ε in both, pyramidal schemes are or may be used.

Dimensionality reduction

From a description point of view, the higher number of features, the better.

Depending on the learning scheme, it may or may not be beneficial. Ideally, one would like to have a feature set that is:

- ⌚ discriminative
- ⌚ decorrelated (no redundancy between features)
- ⌚ invariant to obvious signal modifications
- ⌚ reasonable computational complexity

Features Selection

- ⌚ Brute force subset selection: test all combinations
- ⌚ Single variable classification: sort individual features per classification performance (do not model combined usefulness)
- ⌚ Sequential forward selection: iteratively add the best feature addition to the pool
- ⌚ Sequential backward elimination: iteratively remove the worst feature of the pool

Feature-Space Transformation

The aim is to reduce the number of features while maintaining most of the variability of the data. The most standard way to do this is to perform a Principal Component Analysis (PCA) of the feature set.

ε Let us seek for a transformation matrix T such that:

$$U = T^T \cdot V$$

ε T is a linear projection:

$$T = [c_1, c_2, \dots, c_N]$$

such that

- ε the vectors c_i are the direction of highest variance
- ε the vectors c_i are orthogonal to each others:

$$c_i^T c_j = 0, \forall i \neq j$$



Computation of the PCA

How:

- ① Compute deviations from the mean:

$$B = V - h\mu^T$$

- ② Compute the correlation matrix:

$$C = \frac{1}{N-1} B^* \cdot B$$

- ③ Perform an eigenvalue decomposition:

$$T^{-1} C T = D$$

where D is the diagonal matrix of eigenvalues of C

- ④ sort eigenvectors starting from highest eigenvalue



Computation of the PCA

Usage:

- ε visualization: project data on the 2 or 3 eigenvectors with the highest eigenvalues
- ε features selection: prune out eigenvectors with related eigenvalue lower than 1
- ε features selection: weights of principal eigenvectors are good relevance indicators

Similarity

Task: identify the nearest items of a database of a given seed item

- ε Model features distribution per item
- ε compare models
- ε Validation by considering ranking metrics: k-precision, Mean Average Precision (MAP)



Classification

Task: classify unknown item into pre-defined classes

- ε Feature Extraction
- ε Texture Blocking
- ε Classification decision per block
- ε Majority vote
- ε Validation by computing accuracy using N-folds cross validation: train on N-1 folds, test on 1 fold

Classification schemes

Generative classification scheme: model the joint probability distribution $p(v, l)$ of v the features and l the labels

- ε Gaussian Mixture Models
- ε Hidden Markow Models
- ε ...

Discriminative classification scheme: model the separation

- ε K-Nearest Neighbors
- ε Support Vector Machines (SVM)s
- ε ...