# On the use of statistical tools for speech and musical audio processing

## Mathieu Lagrange

**Analyse / Synthèse Team, IRCAM**

`Mathieu.lagrange@ircam.fr`

**ircam**

**Centre Pompidou**

ATIAM

# Outline

1. **Introduction**
   1. Context and challenges

2. **Past and Present**
   1. Speech
      1. Model
      2. Applications (coding, speaker recognition, speech recognition)
   2. Audio (Music)
      1. Sound models
      2. Retrieving information within songs
      3. Retrieving information across songs

3. **Future**
   1. Polyphony handling
   2. Building and using priors
   3. Joint estimation of several musical parameters

# Outline

1. Introduction
    1. Context and challenges

# Technological Context

« We are drowning in information and starving for knowledge »

R. Roger

- Needs:
  - Measurement
  - Transmission
  - Access

- Aim of a numerical representation:
  - Precision
  - Efficiency
  - Relevance

- Means
  - Mechanical biology
  - Psycho-acoustic
  - Cognition

# Challenges

« Forty-two! yelled Loonquawl. Is that all you've got to show for
seven and a half million years' work?  »

*D. Adams*

Music is a great material to study as it is both:

- An object : arrangement de sons et de silences au cours du
  temps

- A function: more or less codified form of expression of :

  o   Individual feelings (mood)

  o   Collective feelings (party, singing, dance)

# Audio Processing: Past and Present

# Vocabulary ?

- STFT

- MFCCs

- Chromas

- K-means

- GMMs

- HMMs

# Outline

1. Introduction
   1. Context and challenges
2. Past and Present
   1. Speech
      1. Model
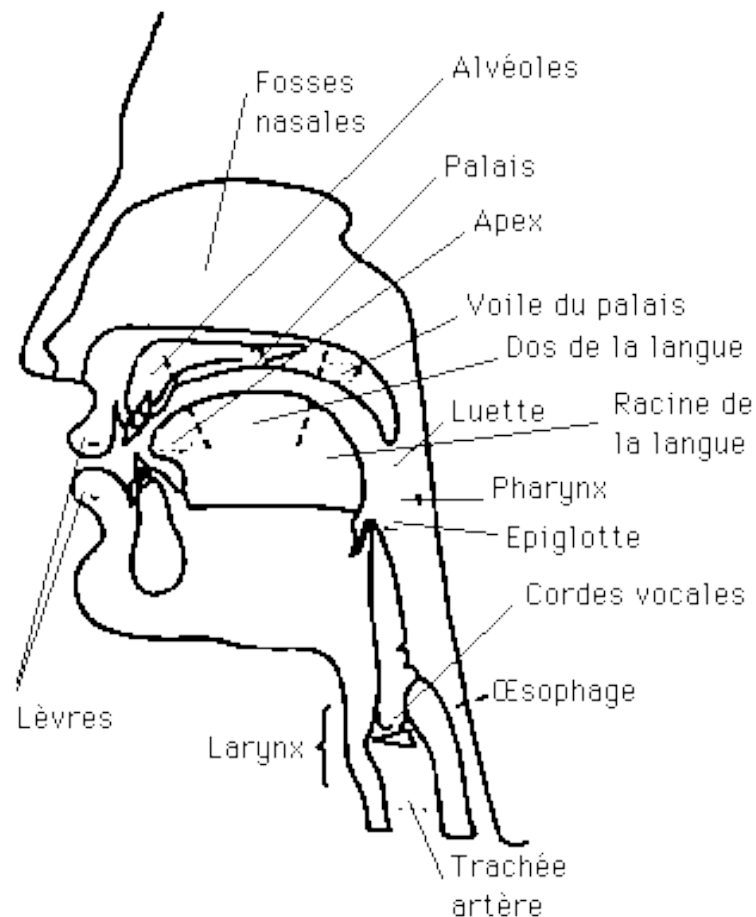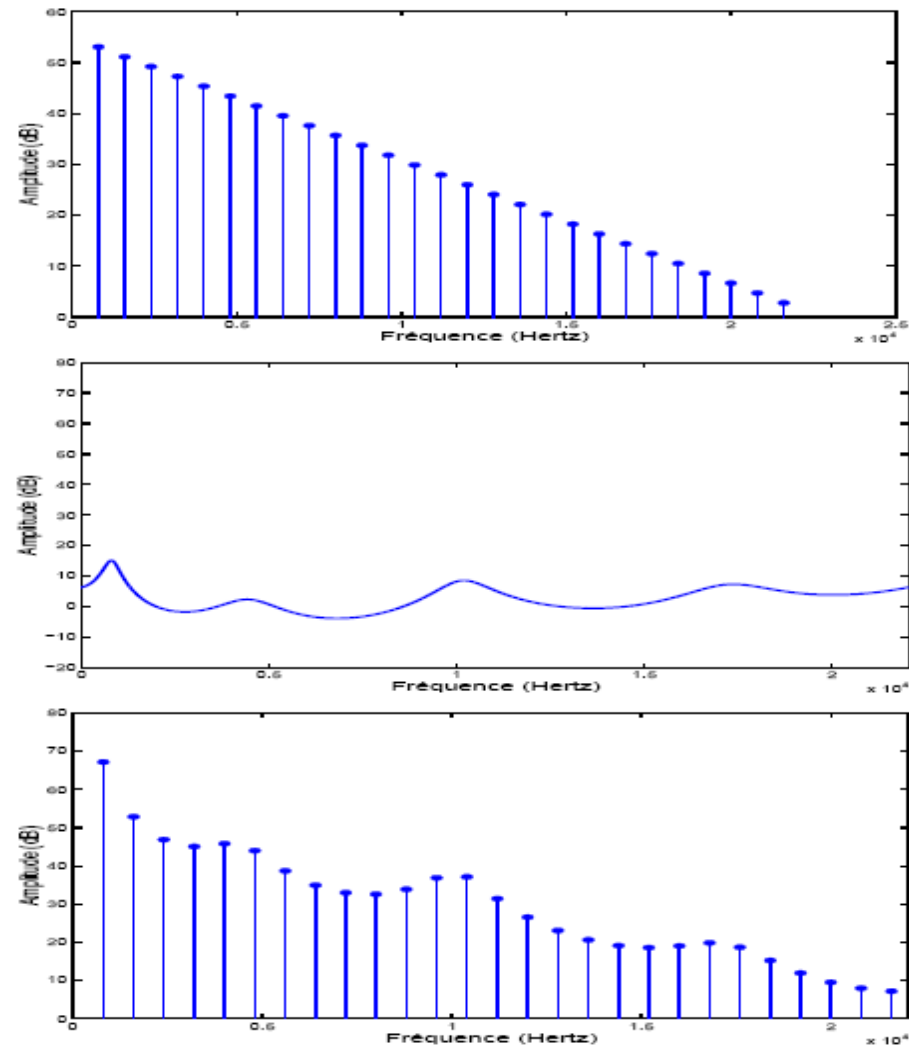      2. Applications (coding, speaker recognition, speech recognition)

# Speech signal

- The speech signal is produced when the air flow coming from the lungs go through the vocal chords and the vocal tract.

    o The size and the shape of the vocal tract as well as the vocal chords excitations are changing relatively slowly

    o The speech signal can therefore be considered as quasi-stationary over short period of about 20 ms.

- Type of speech production

    o Voiced:  <a>, <e>, …

    o Unvoiced:  <s>, <ch>,

    o Plosives: <pe>, <ke>



Fosses nasales
Alvéoles
Palais
Apex
Voile du palais
Dos de la langue
Luette
Racine de la langue
Pharynx
Epiglotte
Cordes vocales
Œsophage
Lèvres
Larynx
Trachée artère

# Source / Filter Model

- In the case of an idealized **voiced** speech signal, the vocal chords are producing a perfectly periodic harmonic signal

- The influence of the vocal tract can be considered as a filtering with a given frequency response whose maximas are called formants.

# Source / Filter Coding

- Algorithm :

    o    Voiced / Unvoiced detection;

    o    Voiced case: the source signal is approximated with a Dirac comb:

        o    a Dirac comb whose successive Diracs are respectively T spaced by T as a spectrum which is a Dirac comb whose successive combs are 1/T spaced.

        o    Parameters : T, gain

    o    Unvoiced: the source signal is approximated by a stochastic signal:

        o    Parameter : gain.

    o    The Source signal is next filtered.

        o    Parameters : filter coefficients.

# « Code-Excited Linear Predictive » (CELP)

For each frame of 20 ms :

o Auto-Regressive coefficients are computed such that the prediction error is minimized over the entire duration of the frame:

$$\widehat{s}[n] = \sum_{i=0}^{d} a[i]s[n-i]$$

$$E = \sum_{n=0}^{N} (s[n] - \widehat{s}[n])^2$$

o Quantified coefficients and an index encoding the error signal are transmitted.

# « Code-Excited Linear Predictive » (CELP)

# Vector Quantization

- Works by dividing a large set of points (vectors in the feature space) into groups

    o    Groups are represented by their centroid

    o    Use of standard k-means algorithms to jointly determine the groups and the centroid

    o    The set of centroids from the codebook

- At the encoding stage

    o    For each residual frame, the closet centroid is determined

    o    The index is transmitted

- At the decoding stage

    o    The centroid is retrieved using the index value

# K-means

- Place k-centroids at random

- Iterate until stabilisation

    o    Determine assignement

    o    Compute centroid

# K-means



(Fig. Wikipedia)

# K-means



(Fig. Wikipedia)

# K-means



(Fig. Wikipedia)

# K-means
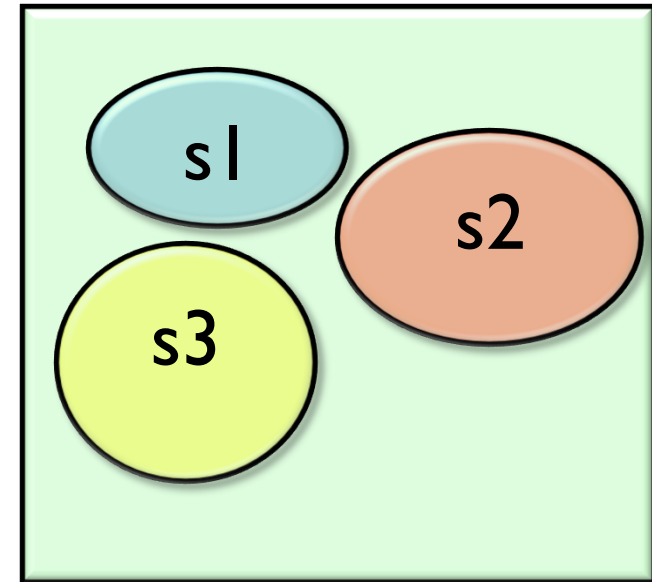
- [Video](Video)



(Fig. Wikipedia)

# K-means

- Place k-centroids at random

- Iterate until stabilisation

    o   Determine assignement
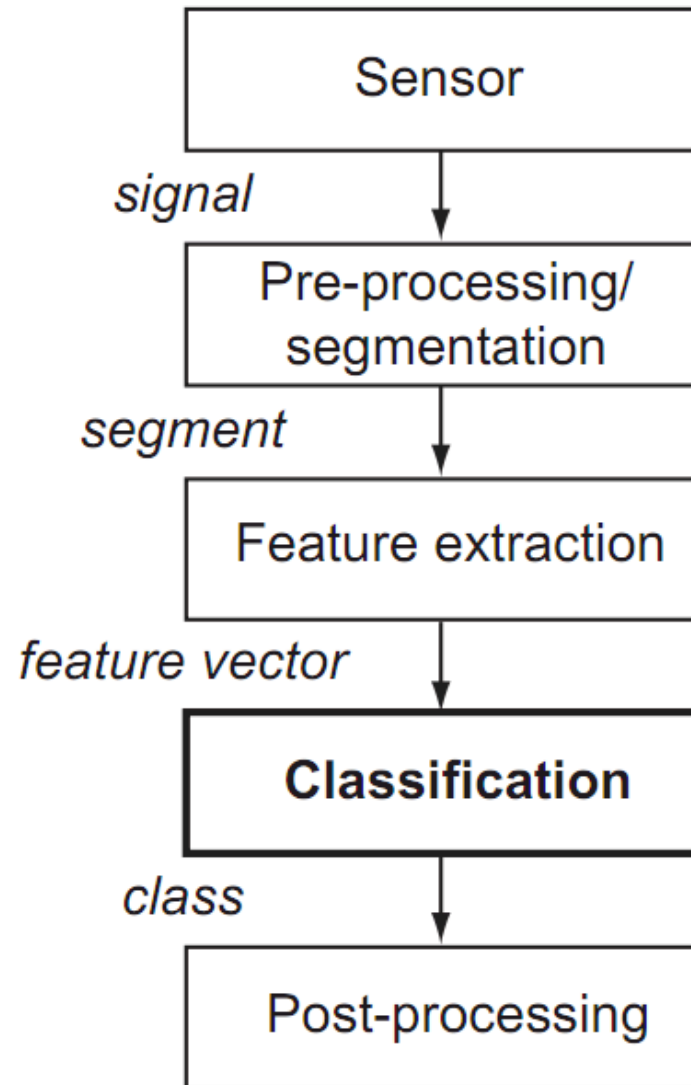
    o   Compute centroid

# Speaker Recognition

- Classical pattern recognition problem

- Specific problems:

  o Open Set / Closed Set: rejection problem

  o Identification / Verification

  o Text Dependency

- Method

  o Feature extraction: model each speech with Mel-Frequency Cepstral Coefficients (MFCCs) and their derivatives.

  o Classification

    o Text independent: Vector Quantization Codebooks or Gaussian Mixture Models (GMMs)

    o Text dependent: Dynamic Time Warping (DTW) or Hidden Markov Model (HMM)

# Classification scheme



Sensor

*signal* → Pre-processing/ segmentation

*segment* → Feature extraction

*feature vector* → **Classification**

*class* → Post-processing

# Features

- The decision system is not usually fed directly with the sound signal

- Infer a reduced set of features

  o Smaller `feature space' (fewer dimensions)

  o Simpler models (fewer parameters)

  o Less training data needed

- Expert knowledge is necessary for efficient inference of meaningful descriptors or features.

  o Meaningful means here that the features

    o Extract from the signal interesting properties for the task at hand

    o Invariance under irrelevant modification

    o have to be nicely handled by the decision system

# Audio Features

- Frequency-dependant information

  o Fourier transform

  

- Mel-Spectrum

  o Robustness to harmonic translation

# Audio Features

- For most classification tasks, we put the focus on the spectral envelope

  o  Speech: formant
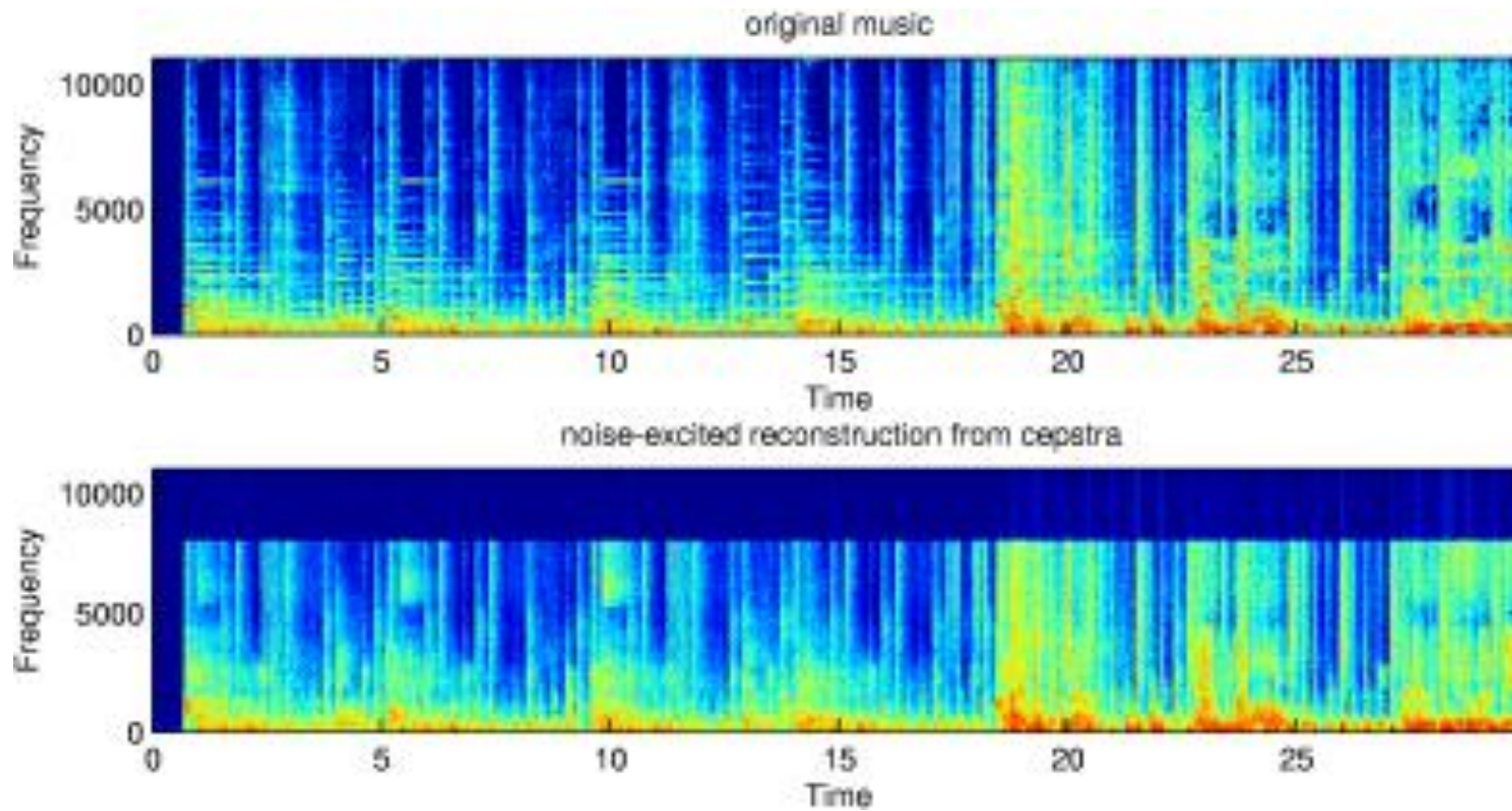
  o  Music: genre

# MFCCs rules ?

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 triangular 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.



Speech signal sampled at 8KHz → Pre - emphasis → Windowing → FFT → Mel scale Filterbank → Take Logarithm → DCT →MFCC CMS →MFCC'

# Example

- Audio



original music

noise-excited reconstruction from cepstra

# Potentials of the DCT step

- Observation of Pols that the main components capture most of the variance using a few smooth basis functions, smoothing away the pitch ripples

- Principal components of vowel spectra on a warped frequency scale aren't so far from the cosine basis functions

- Decorrelates the features.

  o This is important because the MFCC are in most cases modelled by Gaussians with diagonal covariance matrices

# Issues

- The MEL frequency wraping:

  o highly criticized form a perceptual point of view (Greenwood)

  o conceptually: periodicity analysis over data that are not periodic anymore (Camacho)

- The Cepstral Coefficients are COSINE coefficients:

  o cannot shift with speaker size to capture the shift in formant frequencies that occurs as children grow up and their vocal tracts get longer

- Not a sound representation:

  o no way to provide enhancements such as speaker and channel adaptation, background noise suppression, source separation

# Decision System

- From hard assignment to soft assignment

  o  K-means:

# Gaussian Mixture Models (GMMs)

- The data is modeled as a weighted sum of Gaussians

$$g(\boldsymbol{x}, \boldsymbol{\Phi}) = \sum_{k=1}^{g} \pi_k f(\boldsymbol{x}, \boldsymbol{\theta}_k),$$

- Estimation of the weights, means and variances of the Gaussians can be done by maximizing the log-likelihood

$$L(\mathbf{x}; \boldsymbol{\Phi}) = \sum_{i=1}^{n} \log \left( \sum_{k=1}^{g} \pi_k f(\boldsymbol{x}_i, \boldsymbol{\theta}_k) \right).$$

- Usually done with the E-M algorithm

  o   E-step : expectation

  o   M-step : maximisation

# E-M example

- Start



(Fig. From A. Moore's Tutorial)

# E-M example

- 1-st iteration

# E-M example

- 2-nd iteration



(Fig. From A. Moore's Tutorial)

# E-M example

- 3-rd iteration



(Fig. From A. Moore's Tutorial)

# E-M example

- 4-th iteration
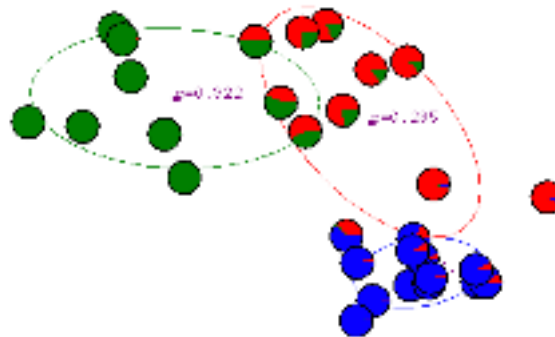


(Fig. From A. Moore's Tutorial)

# E-M example

- 5-th iteration
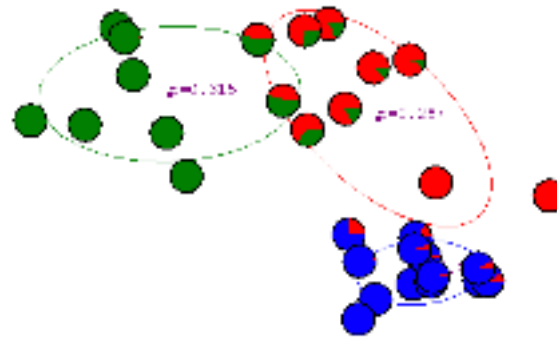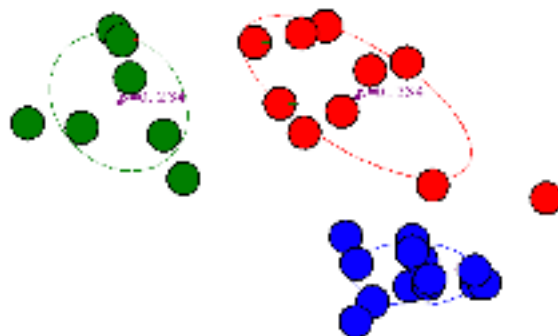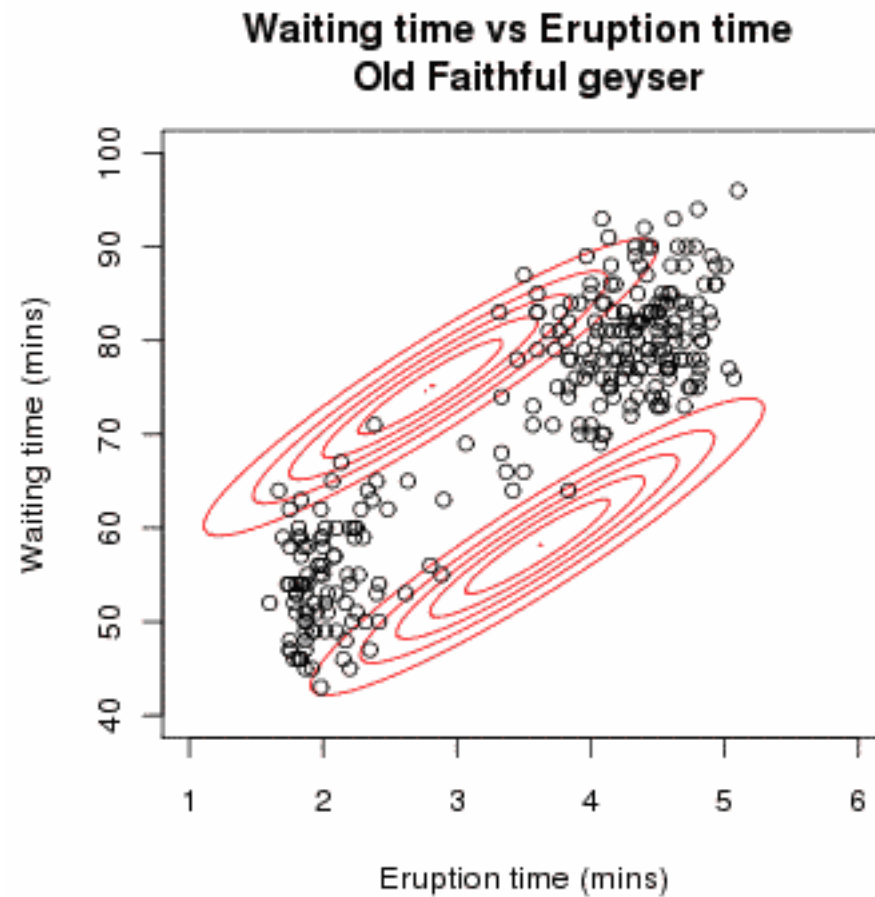
# E-M example

- 6-th iteration

# E-M example

- 20-th iteration

# Density Estimation



Waiting time vs Eruption time
Old Faithful geyser

(Fig. Wikipedia)

# GMMs for the speaker recognition task

- Given a density of probability estimated for each speaker

  o  Search the one that best explains the observed features

- Recent systems are more complex than this

  o  Universal Background Model (UBM)

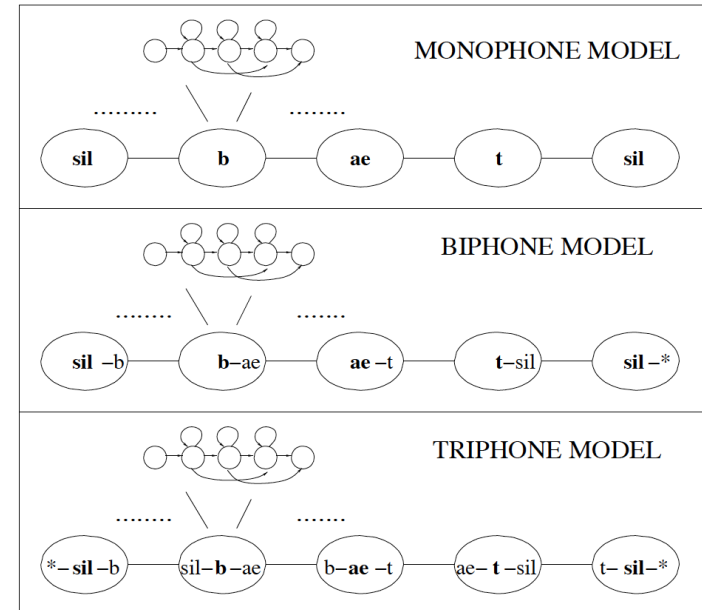  o  Nuisance Attribute Projection (NAP)

# Speech recognition

- The aim of an Automatic Speech Recognizer (ASR) is to

  o  Output the spoken words

  o  Using the speech signal only.



(Fig. from HTK documentation)

# Speech recognition

- An Automatic Speech Recognition System is typically decomposed into:

  o Feature Extraction: MFCCs

  o Acoustic Models: HMMs trained for set of phones

  o Each phone is modelled with **3** states

  o Pronunciation dictionary: convert a series of phones into a word

  o Language Model: predict the likelihood of specific words occurring one after another with n-grams



(Fig. from HTK documentation)

# Summary

- A convenient way of modeling sound is to split the sampled signal into overlapping frames within which the signal is considered as stationary

- Speech encoding:

  o reduced set of parameters that is necessary to synthesize a perceptually similar signal

  o Non parametric: Vector Quantization (K-means)

- Speaker recognition

  o Need to abstract the signal into a meaningful set of parameters: MFCCs

- Speech recognition

  o Sequentiality is important: from GMMs to HMMs

# Outline

1. **Introduction**

    1. Context and challenges

2. **Past and Present**

    1. Speech

        1. Model

        2. Applications (coding, speaker recognition, speech recognition)

    2. Audio (Music)

        1. Sound models

        2. Retrieving information within songs
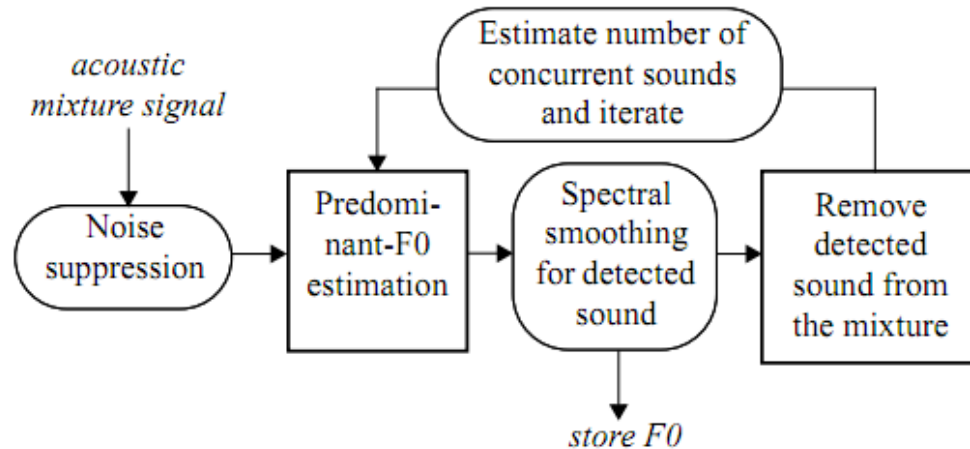
        3. Retrieving information across songs
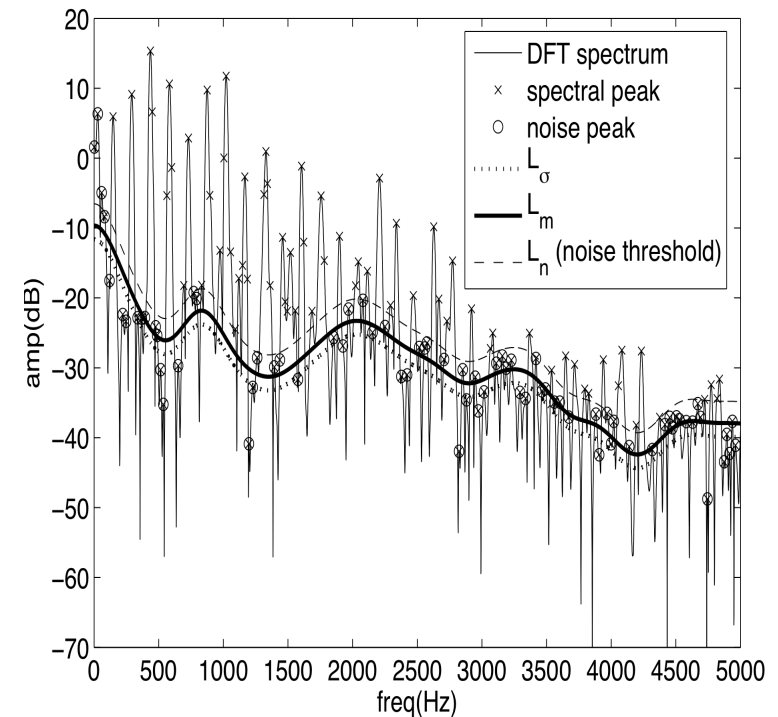
# Retrieving information within songs

- Harmony

  - Pitch Tracking

  - Melody estimation

  - Multi-F0 estimation

  - Chord estimation

- Rhythm

  - Onset detection

  - Tempo estimation

  - Beat tracking

- Orchestration

  - Instrument recognition

# Multi-F0 Estimation

- From an observed spectrum, we want to estimate the fundamental frequency (f0) of each note.

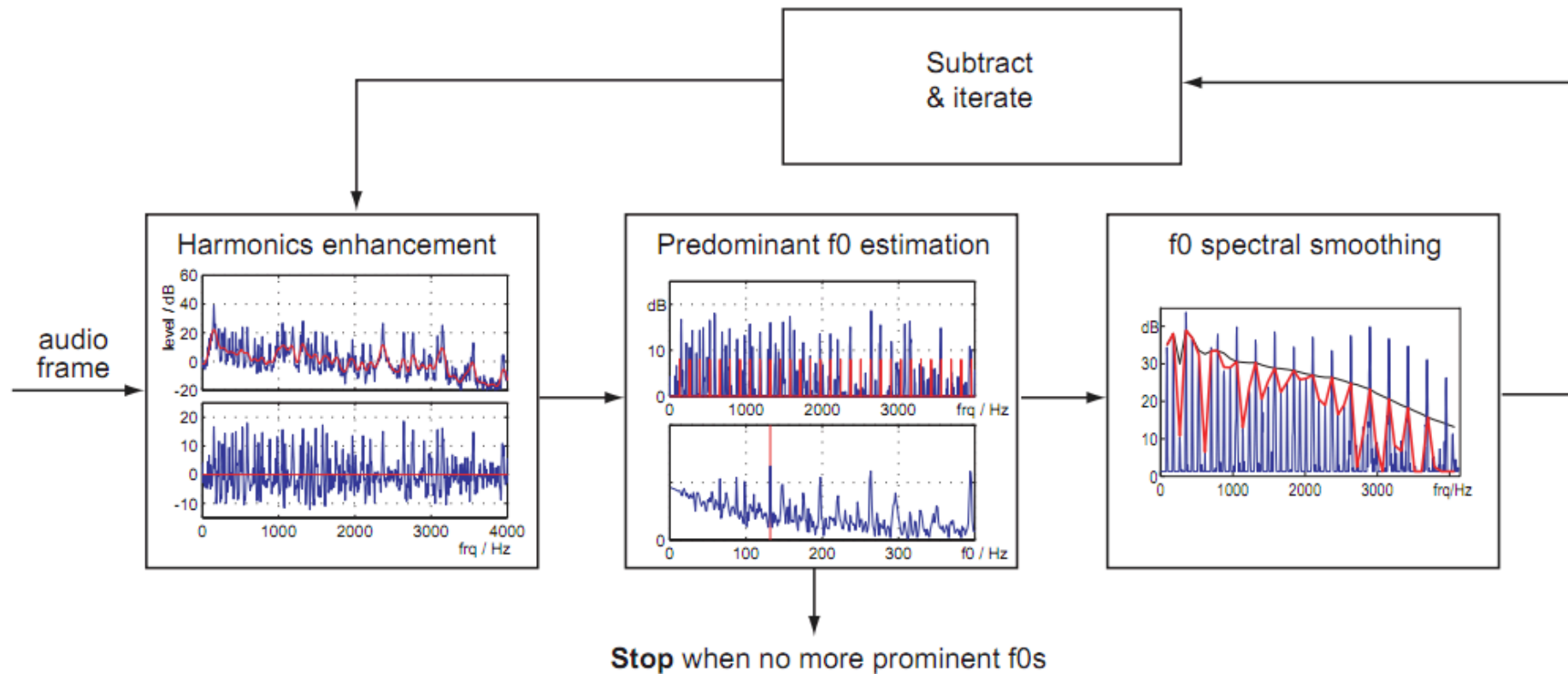    o Most algorithms perform an iterative search:

    o Estimate the dominant f0

    o Remove its contribution



(Fig. from Klapuri 2004)

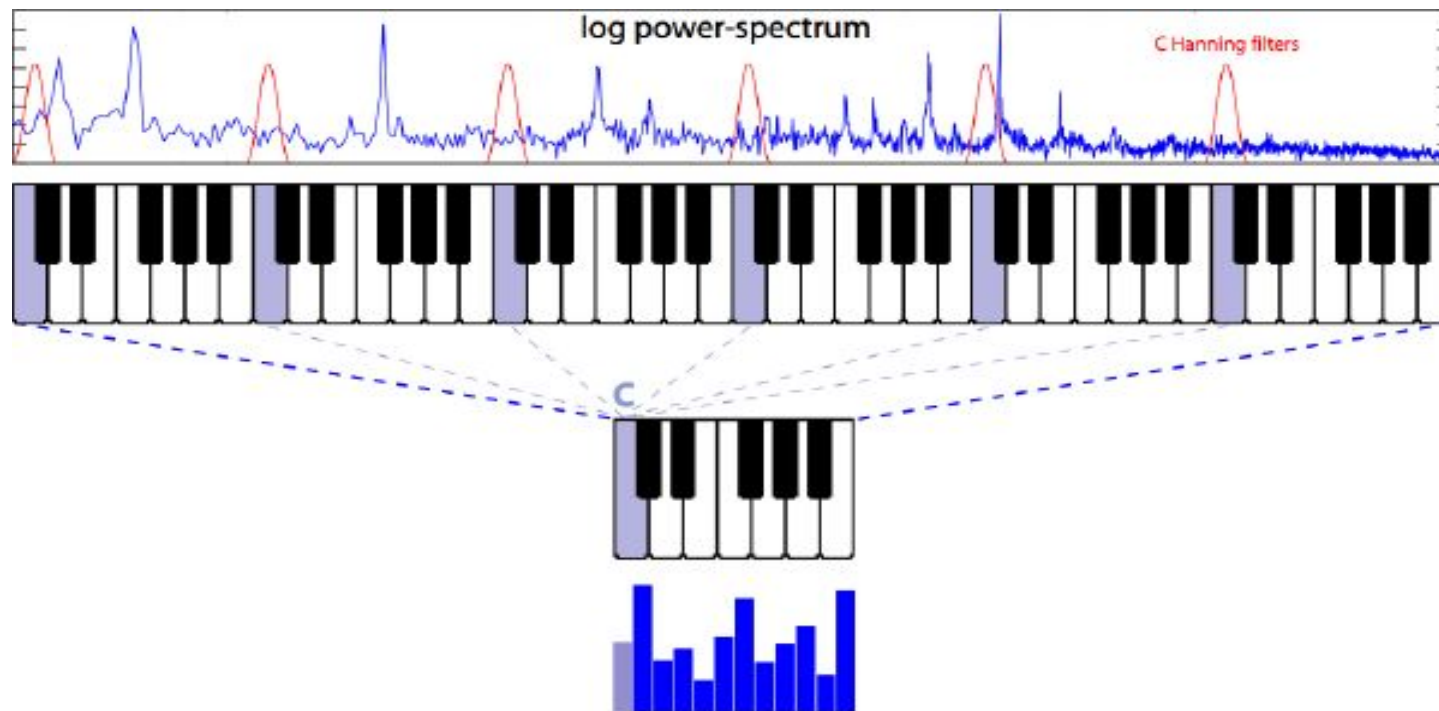(Fig. from Yeh 2010)

# Multi-F0 Estimation

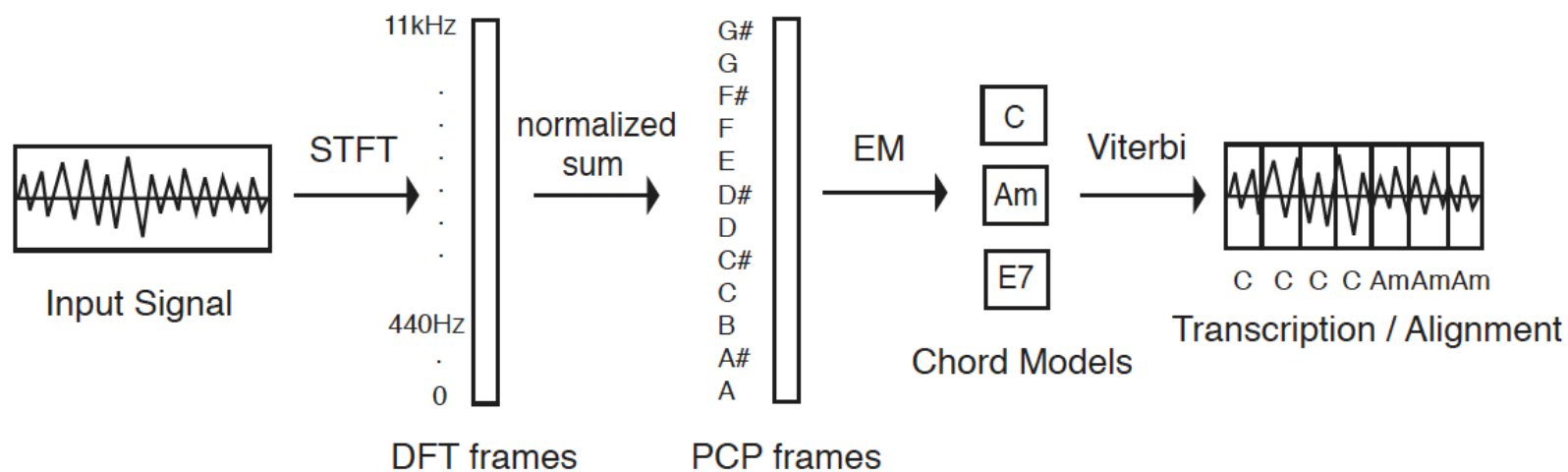# Multi-F0 Estimation



(Fig. from Yeh 2010)

# Chord Estimation

- Relies on a Chroma representation


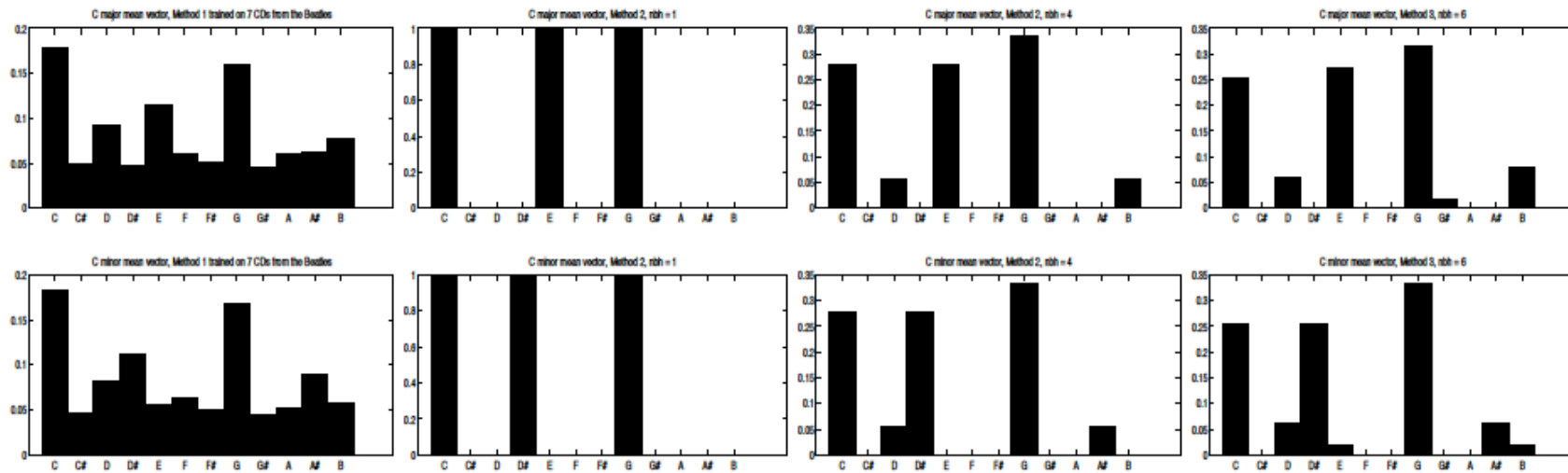
(Fig. from Padapopoulos 2007)

# Chord Estimation

- Matches observed chromas to chords templates

# Issues with template models



- World is dirty, so are our models
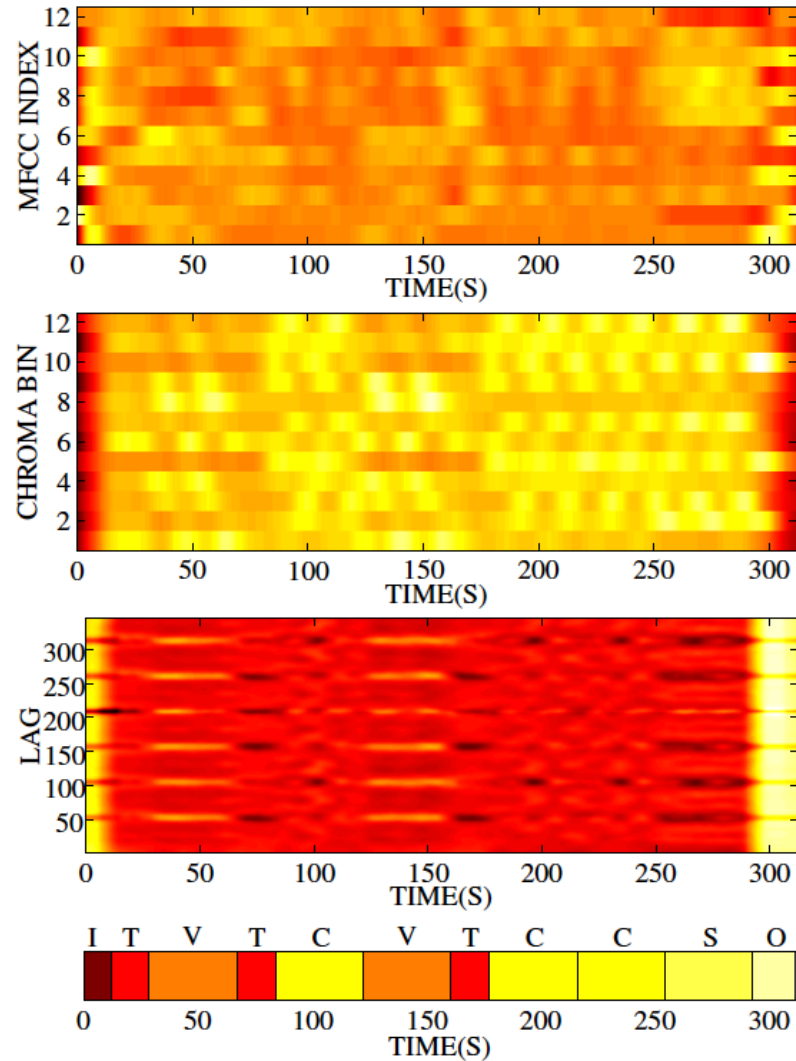
- Our models are clean, let us clean the world

# Structure Analysis

- Aims at estimating the musical structure

  o Ex.: Intro, Verse, Chorus, Verse, Chorus, Chorus, Outro

- Method:

  o Compute features

  o Compute the similarity between every features

  o Perform segmentation based on

    o Novelty

    o Homogeneity

    o Repetition

Paulus'10     Pauls J., Muller M. and Klapuri A.. AUDIO-BASED MUSIC STRUCTURE ANALYSIS
11th International Society for Music Information Retrieval Conference (ISMIR 2010)

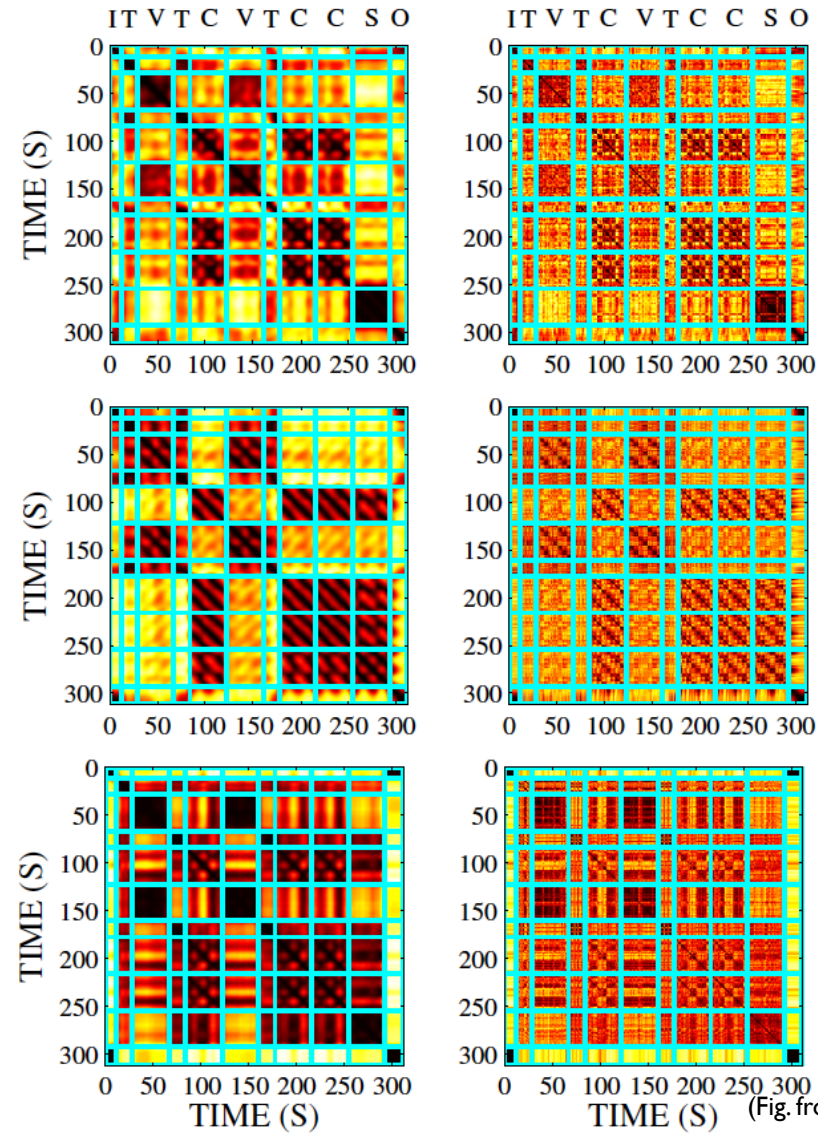# Structure Analysis

- Feature extraction

  o Timbre: MFCCs

  o Harmony: Chromas

  o Rhythm.

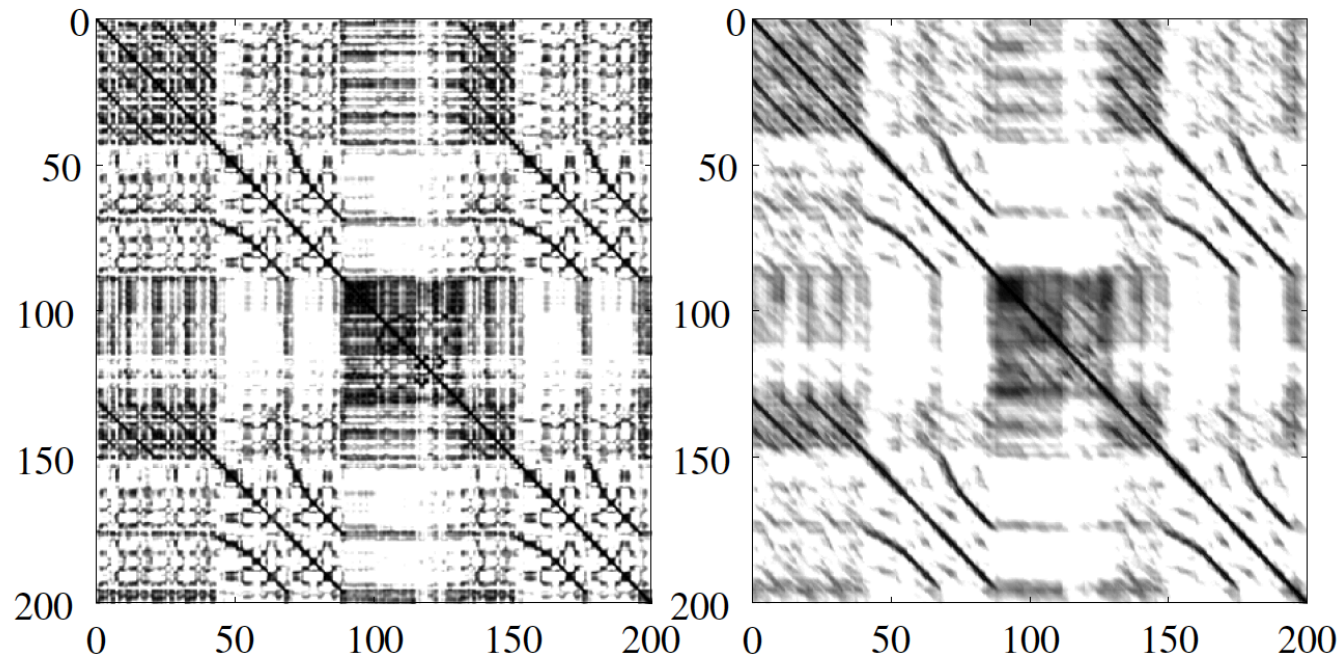

(Fig. from Paulus 2010)

# Structure Analysis

- Self-Similarity Matrix

  o For the 3 features

  o At different granularity



(Fig. from Paulus 2010)

# Structure Analysis
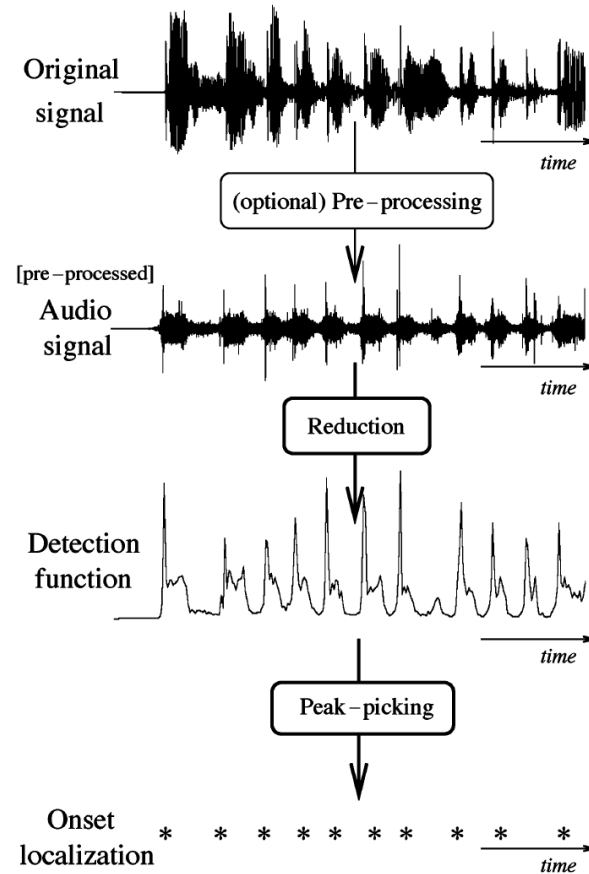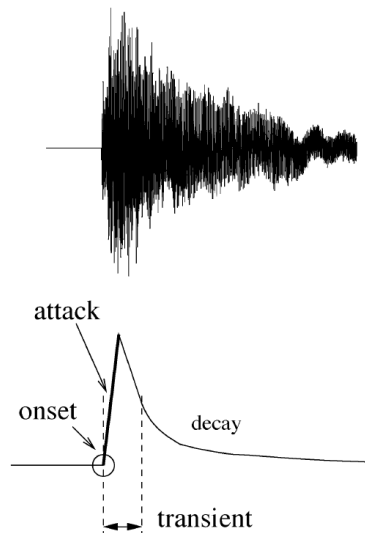


(Fig. from Paulus 2010)

- segmentation based on

  o   Novelty

  o   Homogeneity

  o   Repetition

# Onset Detection

- What is an onset ?

- How do we extract it ?

Bello & al

J. P. Bello, L. Daudet, S. Abdallah, C. Duxbury, M. Davies, and M. B. Sandler, "A tutorial on onset detection in music signals," *Speech and Audio Processing, IEEE Transactions on, vol. 13, no. 5, pp. 1035–1047, Sep. 2005.*
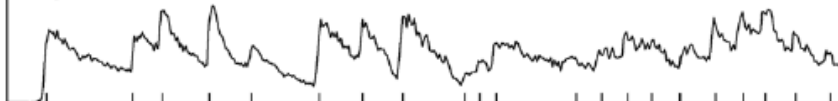
# Onset detection

# Tempo estimation



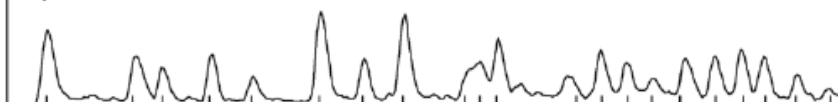**Onset Strength Envelope (part)**

**Raw Autocorrelation**

**Windowed Autocorrelation**

Secondary Tempo Period

Primary Tempo Period

# Beat tracking

- From an estimate of the tempo

  o   Infer the beat position

  o   Possibly the down beat

# Retrieving information across songs

- Classification

  o   Genre recognition

  o   Tag inference

- Similarity

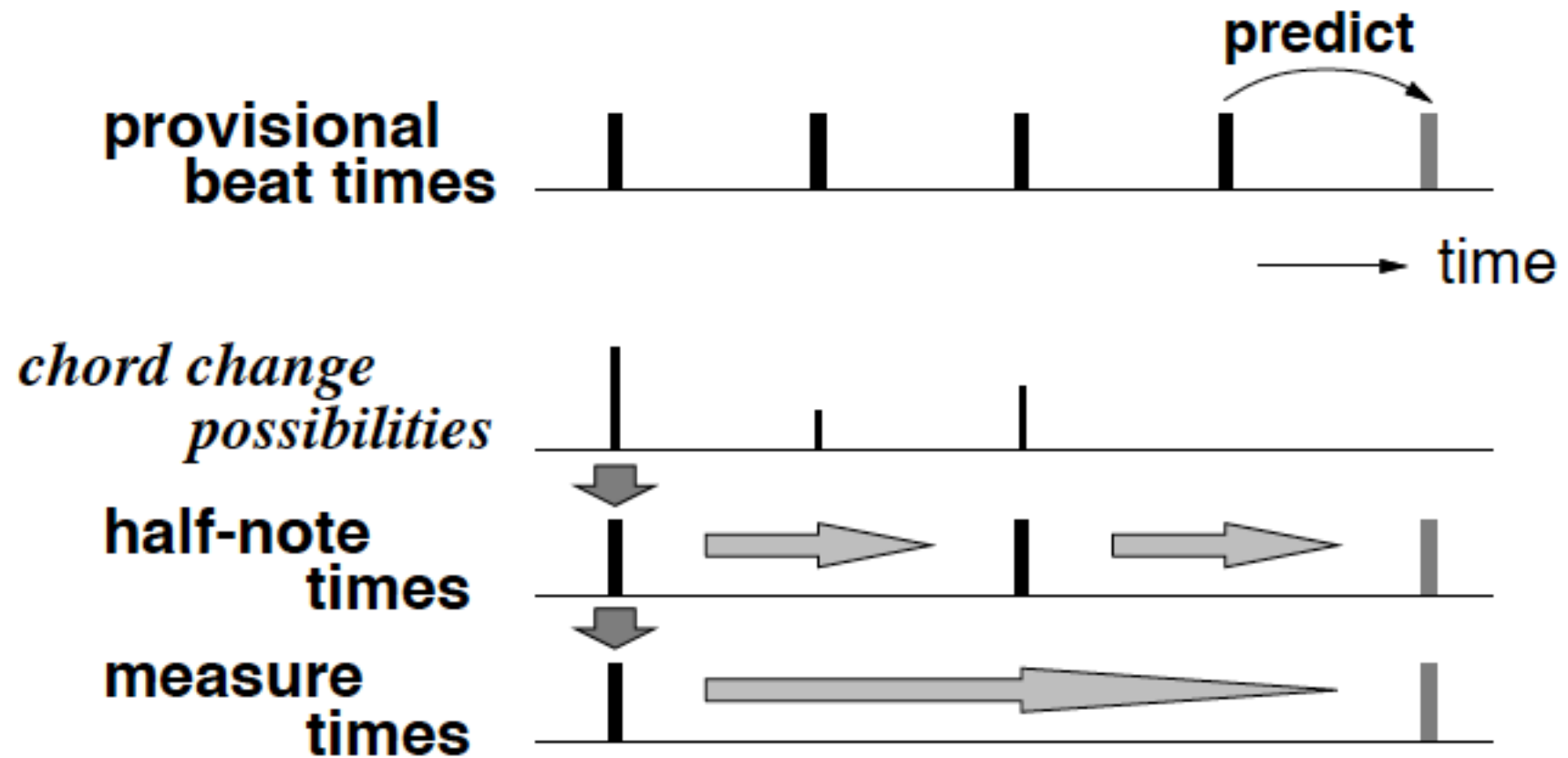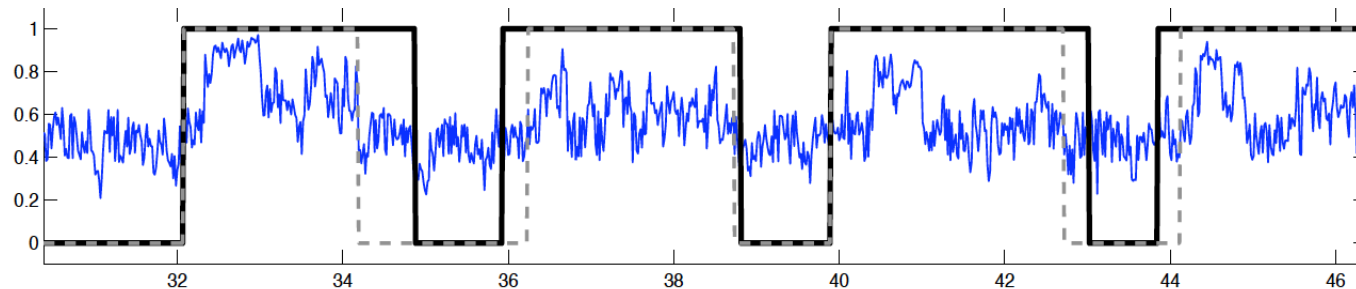  o   Music similarity

  o   Cover detection

# Classification

- Method: [Tzanetakis'02]

  o Agree on mutually exclusive set of tags (the ontology)

  o Extract features from audio (MFCCs and variations)

  o Train statistical models:

    o Due to the high dimensionality of the feature vectors discriminatives approaches are prefered (SVMs)

- Segmentation

  o Smoothing decision using dynamic programming (DP)



(Fig. from [Ramona07])

Tzanetakis'02

Tzanetakis, G. Cook, P. Musical Genre Classification of Audio Signals
IEEE TRANSACTIONS ON SPEECH AND AUDIO PROCESSING 2002

# From parametric to non-parametric
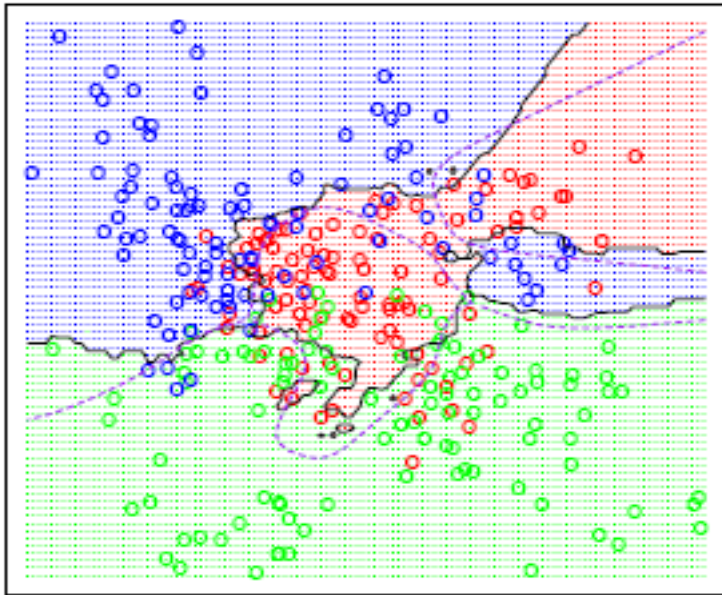
- K-Nearest Neighbors: simple but effective non parametric approach to classification

- Assuming given

  o   A metric that defines the similarity of 2 items

  o   Some labeled items



k=15



k=1

# Issues with k-NN

- Do not scale to large problems (lots of items)

- In high dimensions (lots of features),

  - due to the curse of dimensionality,

  - items tend to be equally far from each others

  - neighbors tend to be non local and meaningless

# Support Vector Machines (SVMs)

- Discriminative approach towards classification

  o In the linear case, SVMs aim at maximizing the distance between margin hyperplanes (dashed lines), called the margin M.



(Fig. from [Ramona Phd])

  o Allows to minimize the structural risk by jointly minimize the

    o Empirical risk

    o Dimension de Vapnik et Chervonenkis

# Kernel-based SVMs

- Data is usually non-linearly separable wich lead to the use of some kernel function to project the data into higher dimensional space

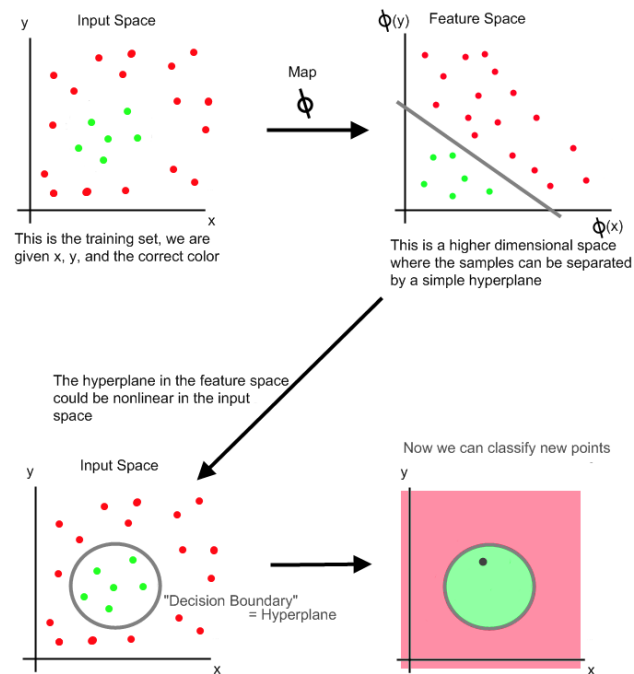If a function $K(.,.)$ is *continous, symmetric* and *positive semi-definite*, then there exists a corresponding mapping $\Phi(.)$ s.t.

$$K(x, y) = \langle \Phi(x), \Phi(y) \rangle \quad \forall x, y \in \mathbb{R}^n$$

# Kernel-based SVMs

Let us consider a polynomial kernel ([Video](Video)) such as

K([x1, x2], [y1, y2]) = x1x2+y1y2+(x1^2+y1^2)(x2^2+y2^-2)

# Multi-Class Discriminative Classification

- Usually performed by combining binary classifiers

- Two approaches:

  o **One-vs-all**: For each class build a classifier for that class versus the rest

  $$f(\mathbf{x}) = \arg\max_i f_i(\mathbf{x}).$$

  o Often very imbalanced classifiers (use asymmetric regularization)

  o **All-vs-all** Build a classifier for each couple of class

  $$f(\mathbf{x}) = \arg\max_i \left( \sum_j f_{ij}(\mathbf{x}) \right).$$

  o A priori a large number of classifiers to build **but** the pairwise classification are faster and the classifications are balanced (easier to find the best regularization)

# Multi-Label Discriminative Classification

- Each object may be tagged using several labels

- Computational approaches

  o Power Sets

  o Binary Relevance (equivalent to one-vs-all)

- Multiple criteria:

  o « Flattening » the ontology

  o Research trend: considering the ontology structure to benefit from co-occurrence labels of different semantic criterion

# Music Similarity

- Question to solve: « Given a seed song, provide us with the entries of the database which are the most similar »

- Annotation type: Artist / Album

- Method:

  o Songs are modeled as a Bag of Features (BoF), usually Gaussian models of MFCCs

  o proximity of GMMs are considered as similiarity measure

    o Diagonal covariance GMMs [Aucouturier'04]:

      - Likelihood (requires access to the MFCCs)

      - Monte carlo sampling

    o Full covariance Gaussian:

      - KL divergence

[Aucouturier'04]    J.-J. Aucouturier and F. Pachet. Improving Timbre Similarity: How High is the Sky? *Journal of Negative Results in Speech and Audio Sciences,* 1 (1), 2004.

# Cover Version Detection

- Question to solve: « Given a seed song, provide us with the entries of the database which are cover versions »

- Annotation: canonical song

- Method: [Serra'08]

  o Songs are modeled as a time series of Chromas

  o Computation of the similarity matrix between the two time series

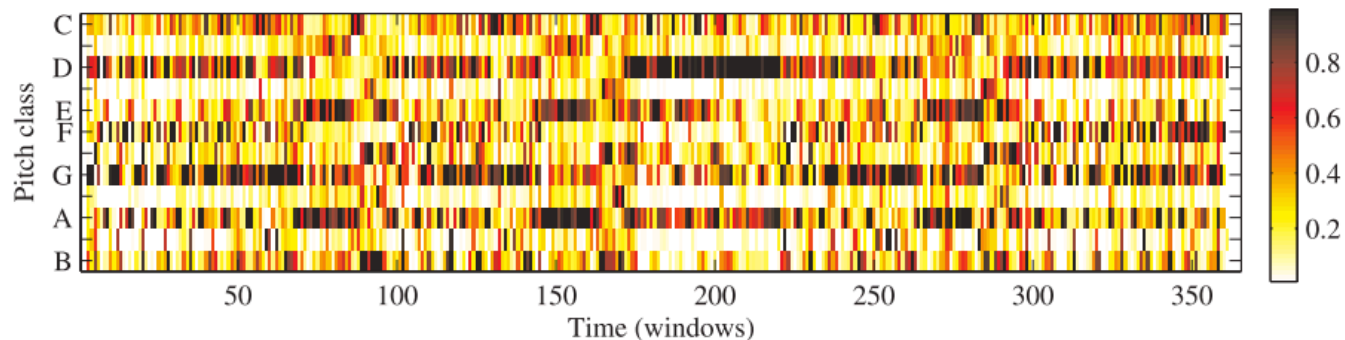  o Similarity is measured using Dynamic Programming Local Alignment
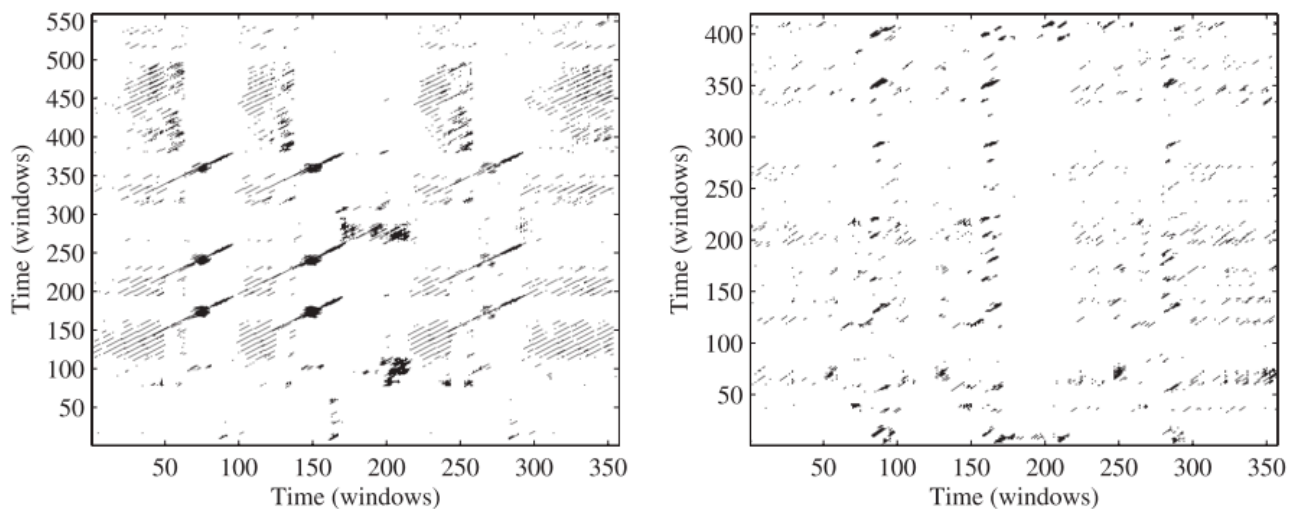
(Fig. from [Serr08])

[Serra'08]    Chroma Binary Similarity and Local AlignmentApplied to Cover Song Identification, IEEE TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING, 2008

# Cover Song Detection

- Chromagram of "Day Tripper"



- Chroma similarity of "Day Tripper" with a cover (left) not a cover (right)

# Future

1. **Issues**
   1. Description of audio and music
      1. Polyphonic
      2. Multiple shapes varying in various ways
   2. Statistical Modeling
      1. Curse of dimensionality
      2. Sense of structure relevant at multiple levels of temporality
2. **Research Trends**
   1. Polyphony handling by Source separation
   2. Building and using priors by performing Auditory Scene Analysis (ASA)
   3. Joint estimation of several musical parameters

# Vocabulary ?

- STFT

- MFCCs

- Chromas

- K-means

- GMMs

- HMMs

# References

- Music information Retrieval (MIR) is an emerging field

  o    Browse ISMIR proceedings

     o    using Google Scholar and use citation index

  o    Stay tuned via Music-IR mailing list

  o    Go to the MIREX webpage to see how things work and are evaluated

# Live coding in Matlab

- You can find the source here:

  o   http://recherche.ircam.fr/equipes/analyse-synthese/lagrange/teaching/atiam11/coursAtiam2011Intro.m

- You will need some external dependencies, web locations are provided in the code

- The code uses cell mode, please look at the Matlab documentation for usage