

M2 ATIAM

# Audio Features Automatic Chord Estimation



Geoffroy Peeters

contact: [geoffroy.peeters@telecom-paris.fr](mailto:geoffroy.peeters@telecom-paris.fr)

Télécom-Paris, IP-Paris, France

## Extracting/ Estimating

### – Automatic Music Transcription

- onset-detection
- pitch, multi-pitch, dominant melody
- chords/guitar tab, key
- rhythm: tempo, meter, beat, downbeat
- musical instrument
- lyrics

### – Auto-tagging

- genre, mood, context, ...

### – Music Recommendation

- similarity, cover-detection

## Processing

### – Source separation, enhancement

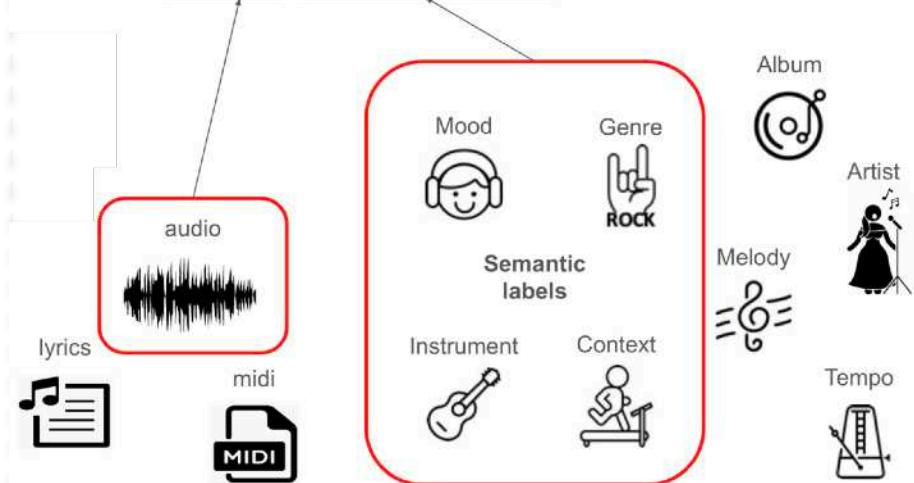
## Generating

### – Audio/ Music generation, transfer

# ISMIR



## Music classification



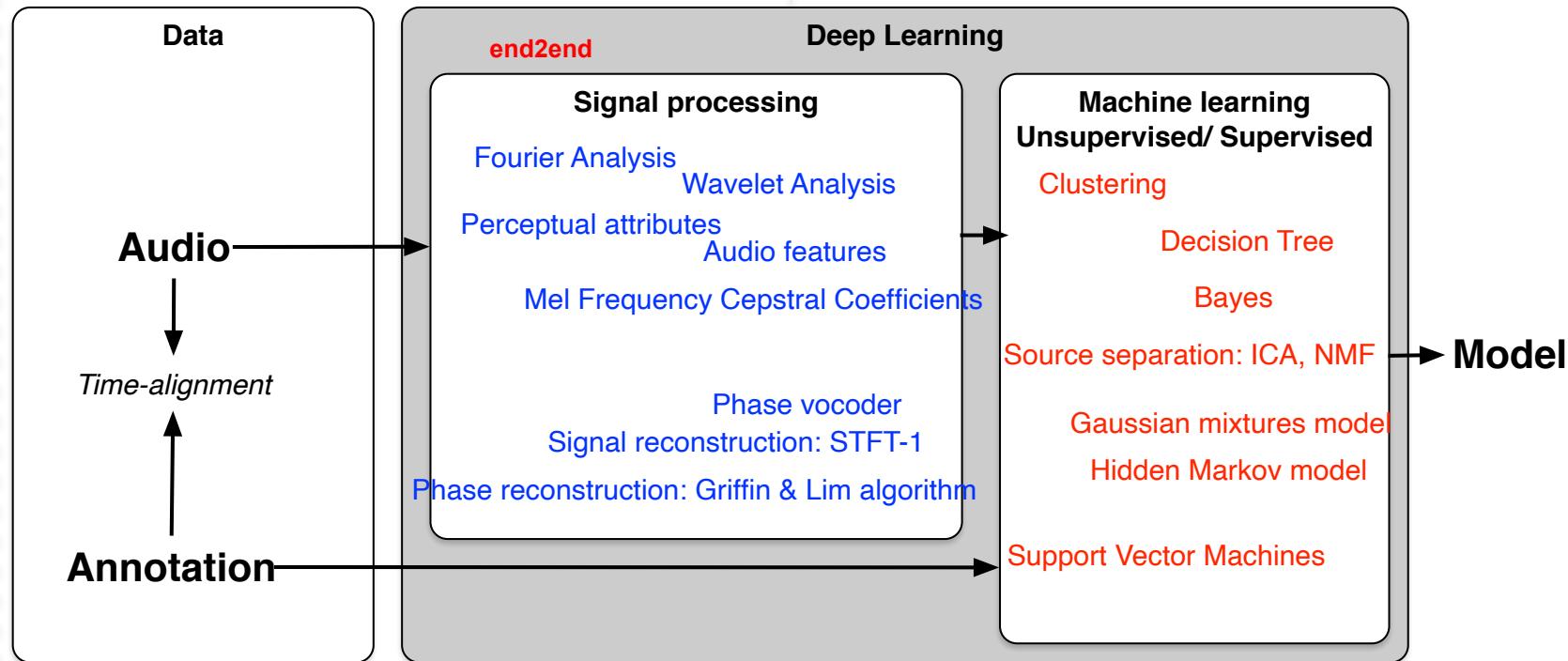
# Various categories of audio content

## Applications

	Speech	Music	Environmental Sounds
Description	Speech to text (ASR) Speaker recognition Speaker diarization	Audio ID (Shazam) Auto-tagging Recommendation Content-description (pitch, chord, tempo) Lyrics alignment/recognition	Acoustic scene classification Sound event localization, detection
Transformation	Speech separation, enhancement Speech coding Speech transformation	Singing separation (Karaoke) Unmixing Style transfer	
Generation	Text to speech	Sound generation Music generation (OpenAI)	Sound texture synthesis
		MIR <a href="http://ismir.net">http://ismir.net</a>	DCASE <a href="http://dcase.community">http://dcase.community</a>

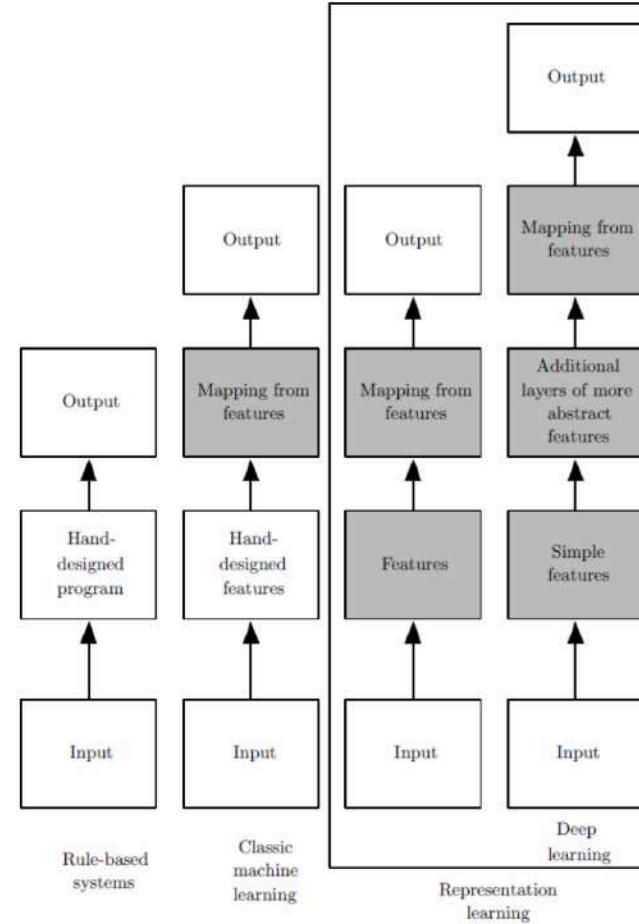


## Traditional approaches for audio



# What is deep learning ?

## Deep learning: learning hierarchical representations



Traditional machine-learning approach

(1) Audio features

# Audio features

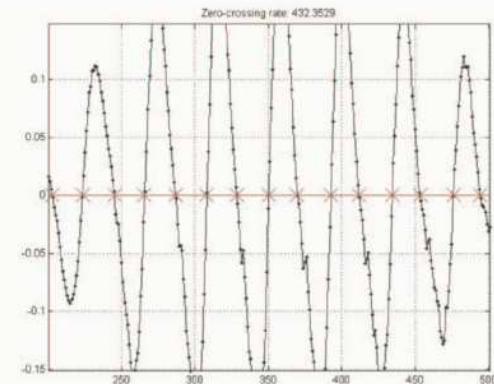
- Various **forms**:
  - **scalar**: spectral centroid, spectral spread, fundamental frequency, spectral roll-off, spectral flux, zero-crossing rate, RMS, ...
  - **vector**: Mel Frequency Cepstral Coefficients, coefficients LPC, coefficients PLP, ...
- Various **time validity**:
  - represent one **frame** of the audio signal → "instantaneous" feature
  - represent the content of a **set of local frames** → texture windows
  - represent **globally** the audio signal
- Highlight different facets of the audio **content**:
  - **timbre** content: Mel Frequency Cepstral Coefficients, LPC coefficients, PLP coefficients, ...
  - **harmonic** content: Pitch Class Profiles/ Chroma, ...
  - **noise** content: Spectral Flatness Measure, ...
  - **rhythmic** content: ...

[G. Peeters. *A large set of audio features for sound description (similarity and classification) in the cuidado project*. Cuidado project report, Ircam, 2004.]

# Audio features examples

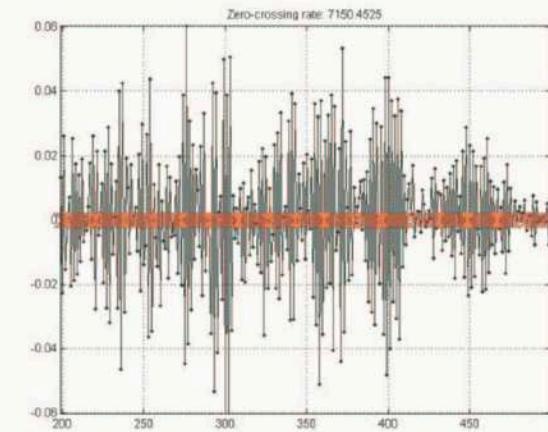
## Zero-crossing rate (zcr)

- Measures the number of times the audio waveform cross the zero-axis
  - $zcr = \frac{1}{N} \sum_{n=1}^N |sign(x_n) - sign(x_{n-1})|$



- Usage: allows to distinguish
  - harmonic sounds → low zcr
  - noise sounds → high zcr

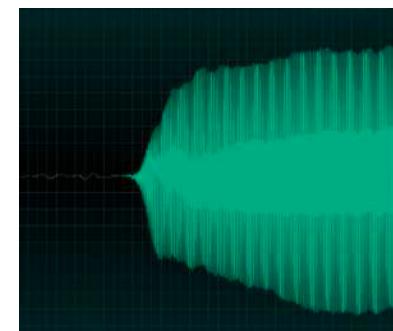
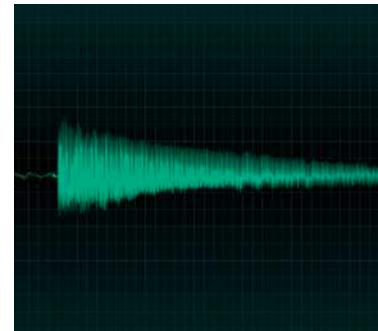
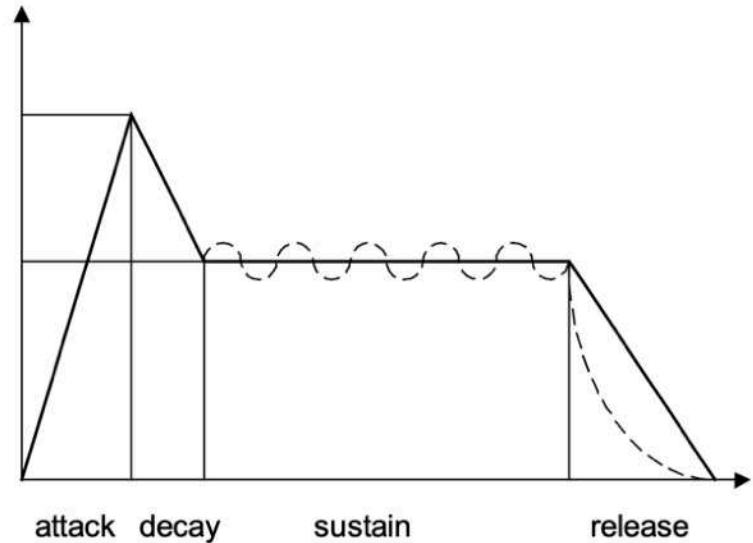
**Figure 12 Zero-crossing rate (=432) during voiced speech region**



**Figure 13 Zero-crossing rate (=7150) during unvoiced speech region**

## ADSR (Attack, Decay, Sustain, Release) temporal enveloppe

- Model of the temporal evolution (enveloppe) of the energy of a musical note
- Usage: allows to distinguish
  - fast attacks (percussive sounds) /slow attacks
  - fast decrease(non-sustained sounds) / slow decrease (sustained sounds)



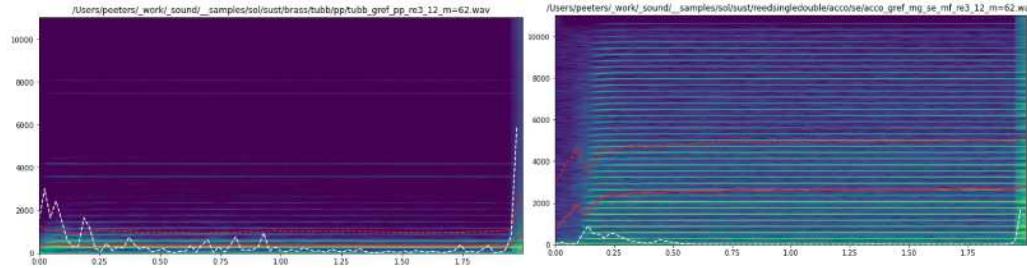
# Audio features examples

## Spectral shape description

### - Spectral centroid

$$\bullet \quad cs = \frac{\sum_k f_k A_k}{\sum_k A_k}$$

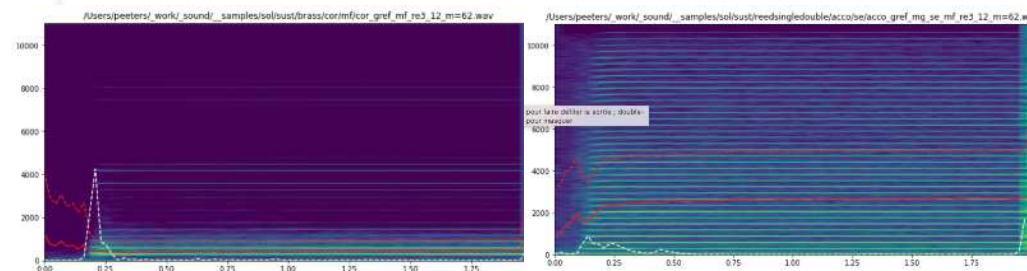
- allows to distinguish between "dull" and "bright" sounds



### - Spectral spread

$$\bullet \quad es = \sqrt{\frac{\sum_k (f_k - cs)^2 A_k}{\sum_k A_k}}$$

- allows to distinguish between "poor" and "rich" sounds

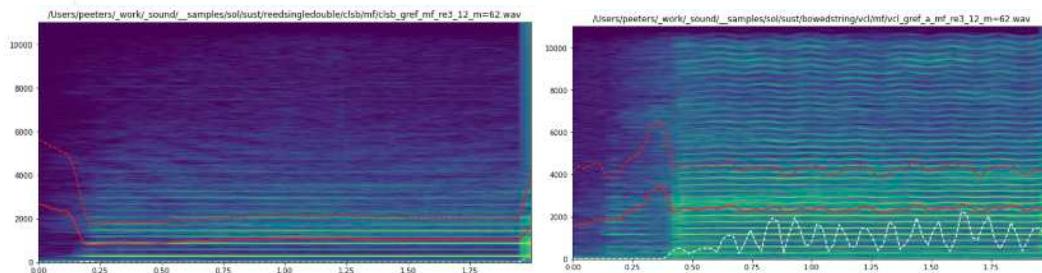


### - Spectral flux

- Measure the temporal variation of the spectrum

$$\underline{fs} = \sum_k (A_k(t) - A_k(t-1))^2$$

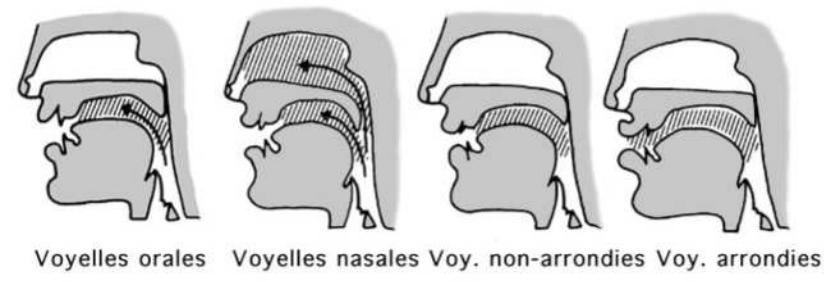
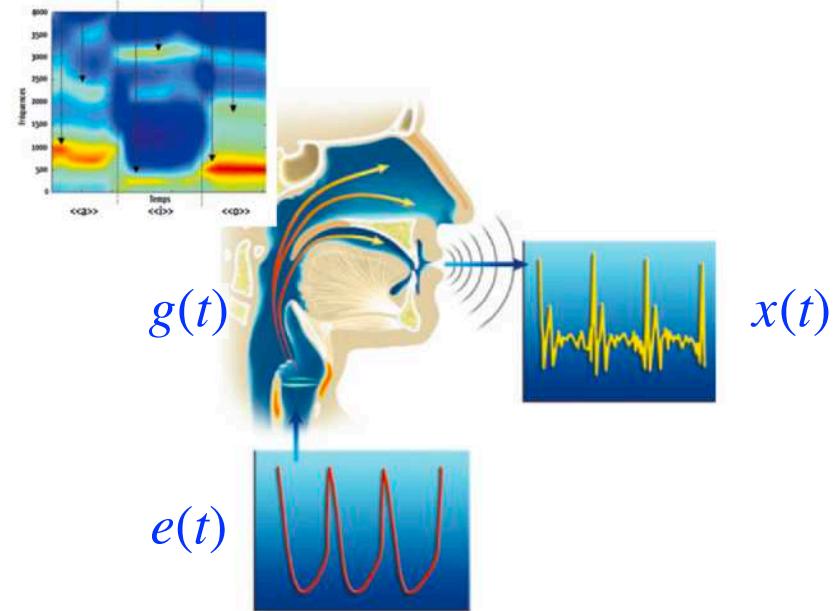
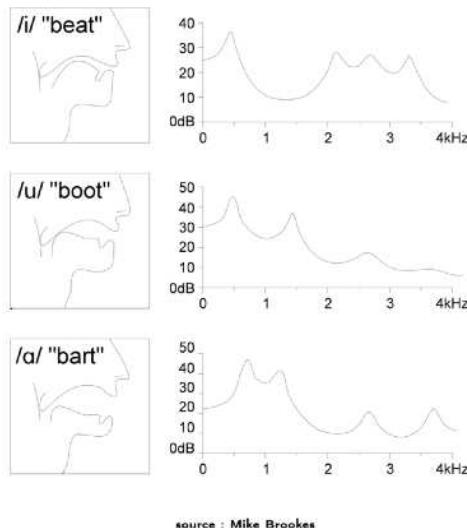
- allows to distinguish between "poor" and "rich" sounds



# Source/Filter model

## Model:

- represent the signal  $x(t)$  as the results of a periodic pulse signal convolved with a resonant filters
- Example:
  - speech (voiced part)
    - vocal folds  $e(t)$
    - mouth (resonance), nose (anti-resonance)  $g(t)$
  - many musical instruments (trumpet)

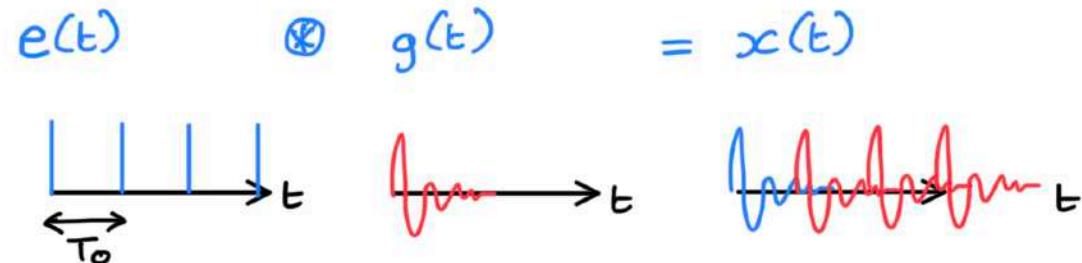


source : [outilsrecherche.over-blog.com](http://outilsrecherche.over-blog.com)

# Source/Filter model

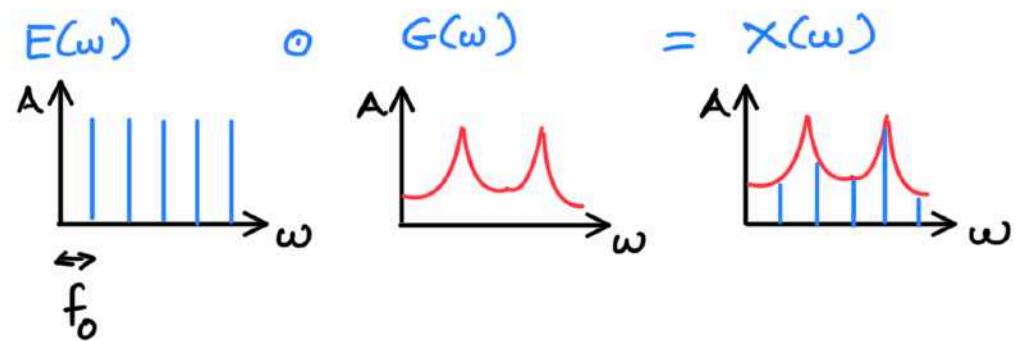
In time:

$$x(t) = e(t) \circledast g(t)$$



In frequency:

$$X(\omega) = E(\omega) \cdot G(\omega)$$



# Audio features examples

## Mel Frequency Cepstral Coefficients (1)

### Complex cepstrum

#### – Goal

- describe the shape of the spectrum (the timbre) of a signal using a reduced set of coefficients

#### – Complex cepstrum $c(\tau)$

$$\begin{aligned} c(\tau) &= TF^{-1} [\log(X(\omega))] \\ &= \frac{1}{2\pi} \int_{\omega} \log[X(\omega)] \cdot e^{j\omega\tau} d\omega \end{aligned}$$

- $\tau$  is named "**quefrency**" (=frequency in reverse order)
- $x(t) \xrightarrow{TF} X(\omega) \xrightarrow{\log} \log(X(\omega)) \xrightarrow{TF^{-1}} c(\tau)$

# Audio features examples

## Mel Frequency Cepstral Coefficients (2)

### Source/filter model

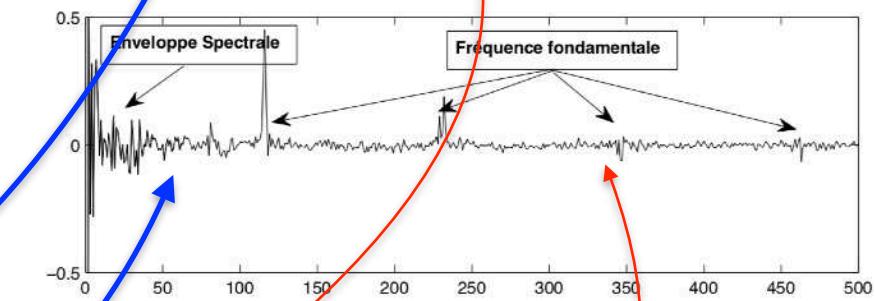
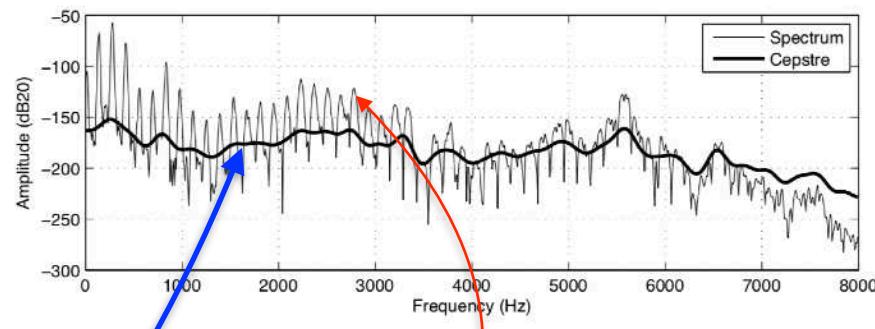
- Source  $e(t)$ : periodic signal
- Filter  $g(t)$ : resonant/ anti-resonant filter

$$x(t) = e(t) \circledast g(t)$$

$$\xrightarrow{TF} X(\omega) = E(\omega) \cdot G(\omega)$$

$$\xrightarrow{\log} \log(X(\omega)) = \underbrace{\log[G(\omega)]}_{\text{slow variations over } \omega} + \underbrace{\log[E(\omega)]}_{\text{fast variations over } \omega}$$

$$\xrightarrow{TF^{-1}} TF^{-1} [\log(X(\omega))] = \underbrace{TF^{-1} [\log[G(\omega)]]}_{\text{energy at quefrency } \tau <<} + \underbrace{TF^{-1} [\log[E(\omega)]]}_{\text{energy at quefrency } \tau >>}$$



# Audio features examples

## Mel Frequency Cepstral Coefficients (3)

### Real cepstrum

- **Real ?** = cepstrum computed on the real part of the log-spectrum

$$X(\omega) = A(\omega) \cdot e^{j\phi(\omega)}$$

$$\log[X(\omega)] = \log[A(\omega)] + j\phi(\omega)$$

$$\Re\{\log[X(\omega)]\} = \log[A(\omega)]$$

$$\text{real cepstrum} = TF^{-1} [\Re\{\log[X(\omega)]\}]$$

$$= TF^{-1} [\log[A(\omega)]]$$

$$c(\tau) = \frac{1}{2\pi} \int_{\omega} \log[A(\omega)] \cdot e^{j\omega\tau} d\omega$$

- The amplitude spectrum  $A(\omega)$  is real and symmetric
  - its Fourier Transform reduces to the real part
    - reduces to the projection of  $\log[A(\omega)]$  on a set of cosinus → Discrete Cosine Transform (DCT)

# Audio features examples

## Mel Frequency Cepstral Coefficients (4)

### Mel Frequency Cepstral Coefficients (MFCCs)

- MFCC ? = real cepstrum computed on the power spectrum  $|X(\omega)|^2$  converted to the Mel scale (a perceptual scale)
- **Why perceptual scales ?**
  - Fourier Transform
    - decomposition on a set of sinusoidal components which frequencies are linearly spaced ( $f_k = 10\text{Hz}, 20\text{Hz}, 30\text{Hz}, \dots \text{Hz}$ )
  - Human hearing:
    - decomposition on a set of filters which frequencies are logarithmically spaced (10, 20, 40, 80, ... Hz).
    - highest resolution in low-frequencies, lowest resolution in high frequencies
    - in speech, formants/resonances are closer together in low frequencies
  - MFCCs allows a more compact representation than the real cepstrum
- **How ?**
  - Use of perceptual scales: Mel-scale, Bark-scale, ERB-filters, Gamma-tone filters
- **Usage ?**
  - MFCCs are the most used features in audio: speech, music, environmental sounds recognition, ...

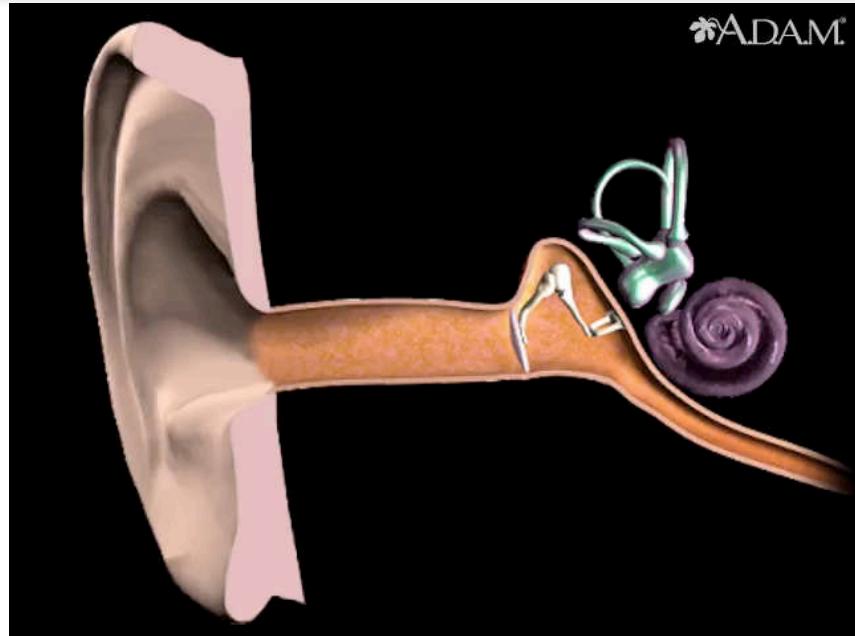
# Audio features examples

## Mel Frequency Cepstral Coefficients (5)

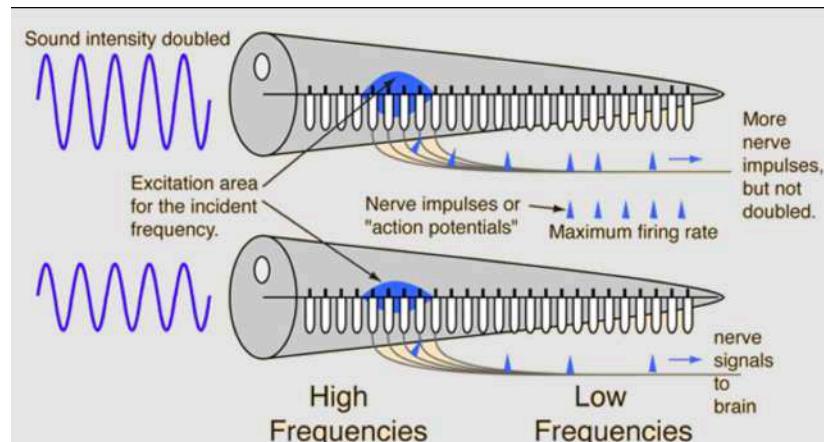
### Human hearing

- Cochlea
- Critical bands
  - perception of two tones at  $f_1$  and  $f_2$
  - perception of a beating-tone at  $\frac{f_1 + f_2}{2}$

$$\cos f_1 + \cos f_2 = 2 \cos \frac{f_1 - f_2}{2} \cos \frac{f_1 + f_2}{2}$$



<https://medlineplus.gov/ency/anatomyvideos/000063.htm>



source: <http://hyperphysics.phy-astr.gsu.edu/hbase/Sound/loud.html>

# Audio features examples

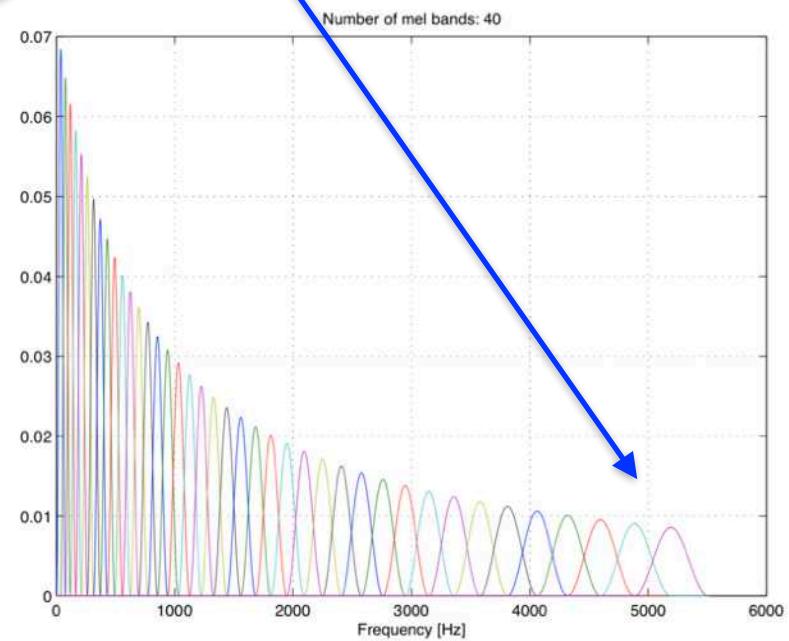
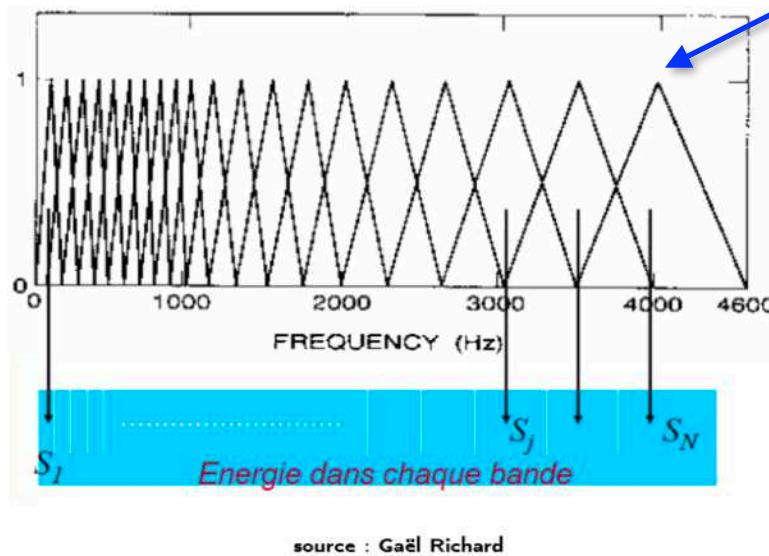
## Mel Frequency Cepstral Coefficients (6)

### Mel scale ?

$$mel(f) = \frac{1000}{\ln 2} \ln \left( 1 + \frac{f}{1000} \right)$$

- Remark: variations of the constant exist

different shapes for the filter: triangular, hanning, tanh



Fant, Gunnar. (1968). *Analysis and synthesis of speech processes*.

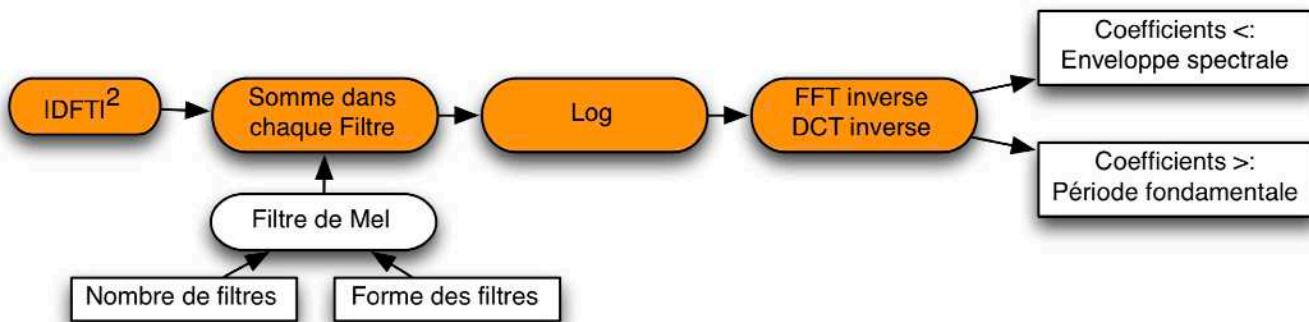
In B. Malmberg (Ed.), *Manual of phonetics* (pp. 173-177). Amsterdam: North-Holland.

# Audio features examples

## Mel Frequency Cepstral Coefficients (7)

### Computation steps for MFCCs

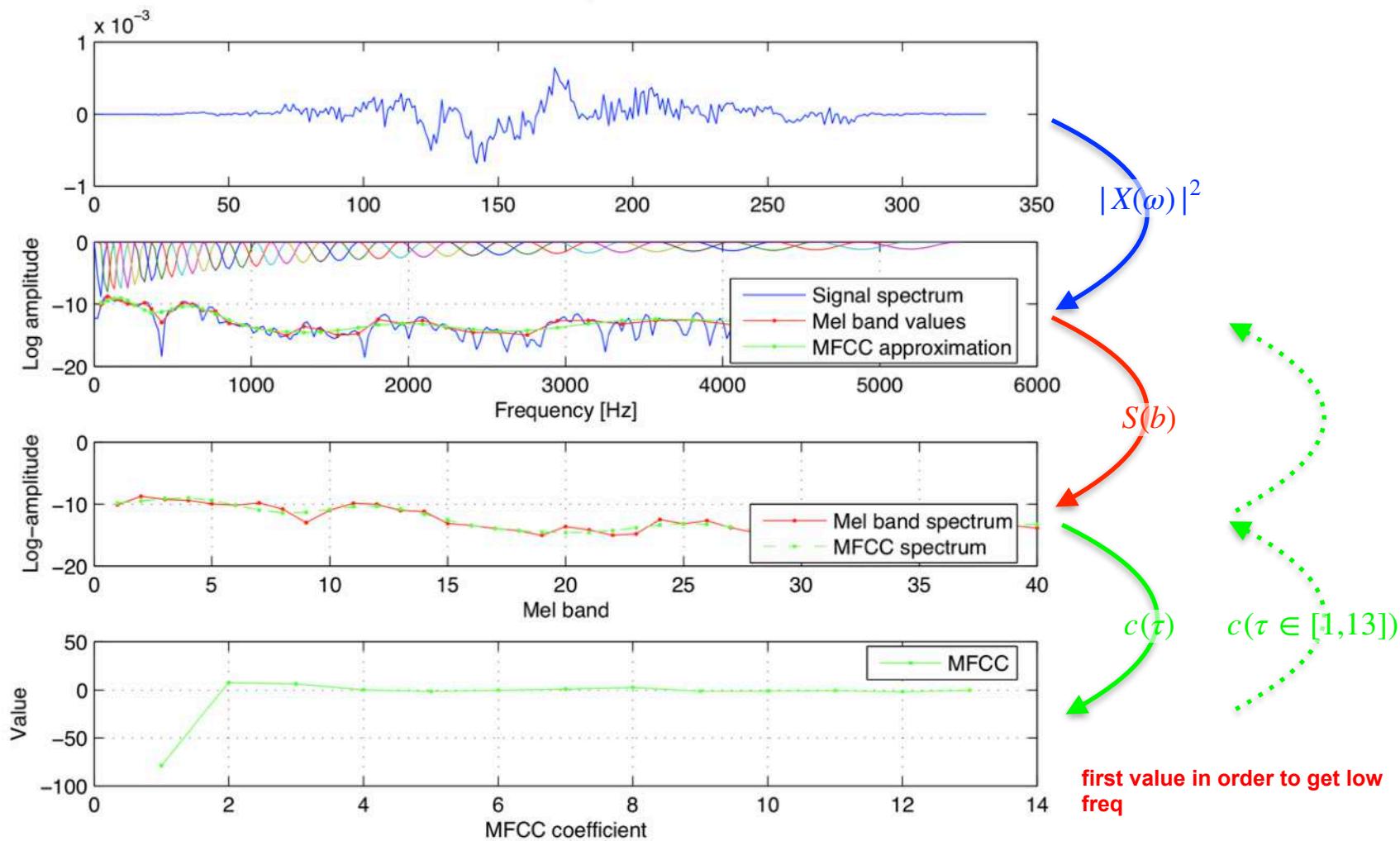
- Compute the power spectrum:  $|X(\omega)|^2$
- Compute the Mel filters:  $H_b(\omega)$  with  $b \in [1, B]$ 
  - choice of the number of filters  $B$ : 40
  - choice of the shape of each filter: triangular, hanning, tanh, ...
- Convert the power spectrum to Mel bands:  $S(b) = \sum_{\omega} |X(\omega)|^2 \cdot H_b(\omega)$
- Convert to logarithmic scale:  $\log(S(b))$
- Compute the IFFT (or the IDCT):  $c(\tau)$
- Select the first coefficients, close to 0 (usually the first 13 coefficients)
  - coefficients close to zero represent the decomposition of the Mel bands content on a set of cosinus with slow variations



# Audio features examples

## Mel Frequency Cepstral Coefficients (8)

### Example of the computation of MFCCs



# Chroma/ Pitch Class Profile (PCP)

# Chroma/ Pitch Class Profile (PCP)

- **Helical model of pitch [Roger Shepard, 1964]**

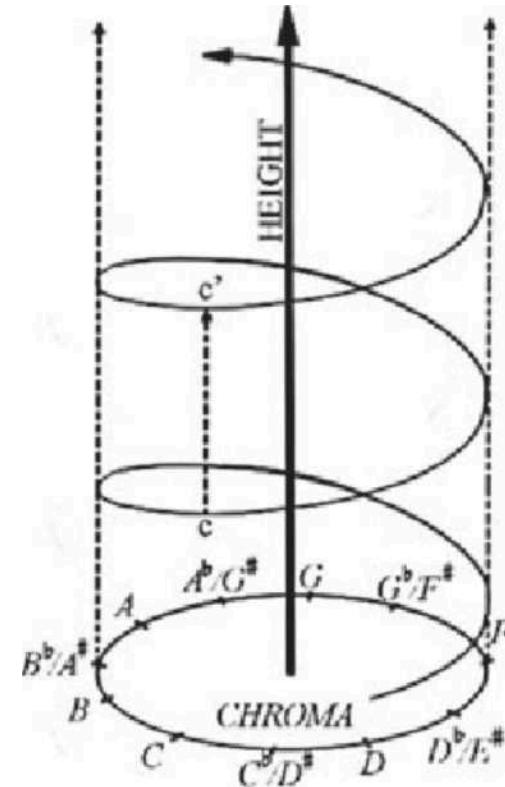
- represents the pitch of a note  $p$  as a two-dimensional structure:
- $p = c + o \cdot 12$ 
  - chroma  $c$  (pitch-class)
  - tonal height  $o$  (octave number)

- **Definition: Chroma - Pitch Class Profile (PCP):**

- represents the harmonic content of the spectrum at time  $n$ ,  $X(k, n)$ , as a vector:
  - $C(c, n) \quad c \in [0, 12[$

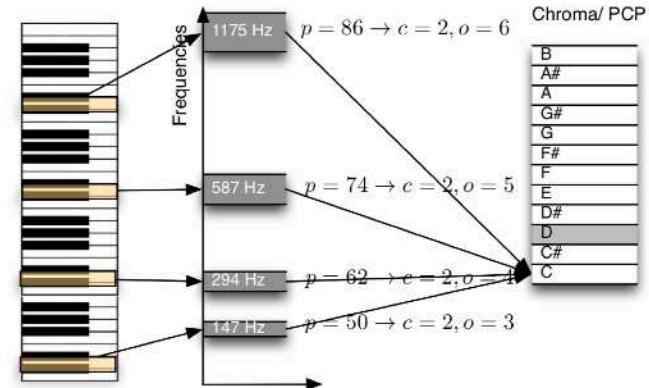
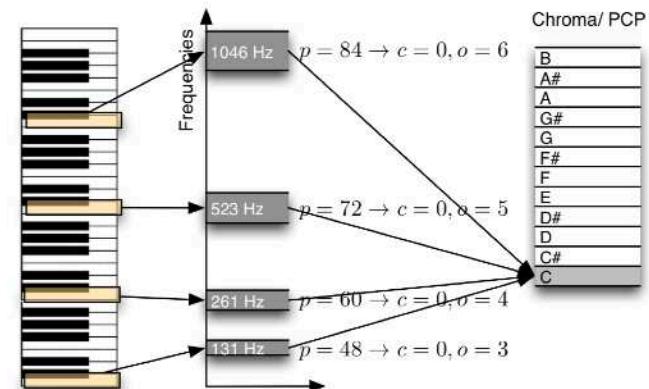
- **Usage:**

- key estimation,
- chord estimation,
- cover detection



# Chroma/ Pitch Class Profile (PCP)

- **Chroma computation  $C(c, n)$** 
  - We sum up the values of the spectrum  $X(k, n)$  for all  $f_k$  which correspond to a given  $c$
  - Relationship between the frequencies  $f_k$  of the DFT and the pitches  $p$  (semi-tone pitches in MIDI-scale)
    - $p(f_k) = 12 \log_2 \left( \frac{f_k}{440} \right) + 69, p \in \mathbb{R}^+$
    - $f(p) = 440 \cdot 2^{\frac{p-69}{12}}$



T. Fujishima. Realtime chord recognition of musical sound: a system using common lisp music. In Proc. of ICMC (International Computer Music Conference), pages 464–467, Beijing, China, 1999.

G. H. Wakefield. Mathematical representation of joint time-chroma distributions. In Proc. of SPIE conference on Advanced Signal Processing Algorithms, Architecture and Implementations, pages 637–645, Denver, Colorado, USA, 1999.

# Chroma/ Pitch Class Profile (PCP)

- Spectral resolution ?**

- Should allow separating adjacent musical notes

– We define the width (at -6 dB) as  $B_w = \frac{C_w}{L_{sec}}$

- If  $f_{min}$  (the lowest frequency we consider in the spectrum) is 50 Hz

- We need to separate  $G\#1$  (51.91Hz) from  $A1$  (55Hz)

$$\rightarrow L_{sec} = \frac{C_w}{B_w} = \frac{2.35}{3.0869\text{Hz}} = 0.7613\text{s}$$

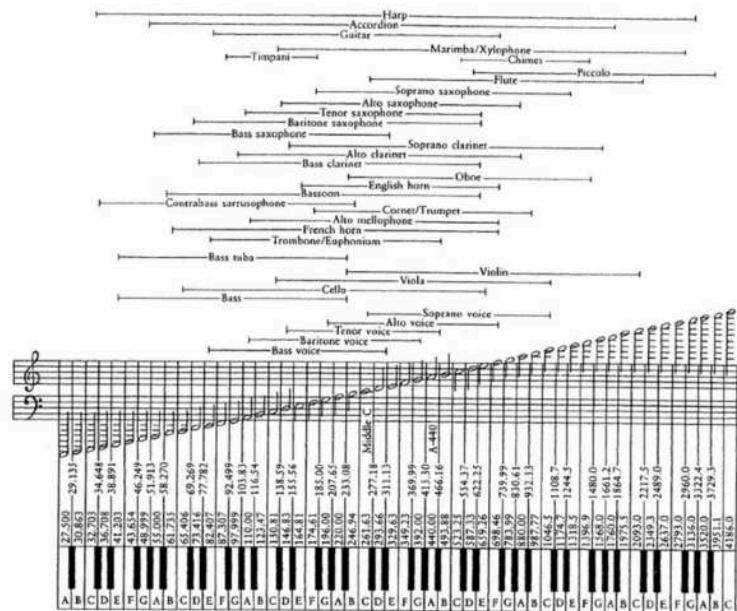
- If  $f_{min}$  is 100 Hz

- We need to separate  $G\#2$  (103.82Hz) from  $A2$  (110Hz)

$$\rightarrow L_{sec} = \frac{C_w}{B_w} = \frac{2.35}{6.1738\text{Hz}} = 0.3806\text{s}$$

- Two possibilities:

- Choice  $L_{sec}$  as a function of  $f_{min}$
- Choose  $f_{min}$  as a function of  $L_{sec}$



# Chroma/ Pitch Class Profile (PCP)

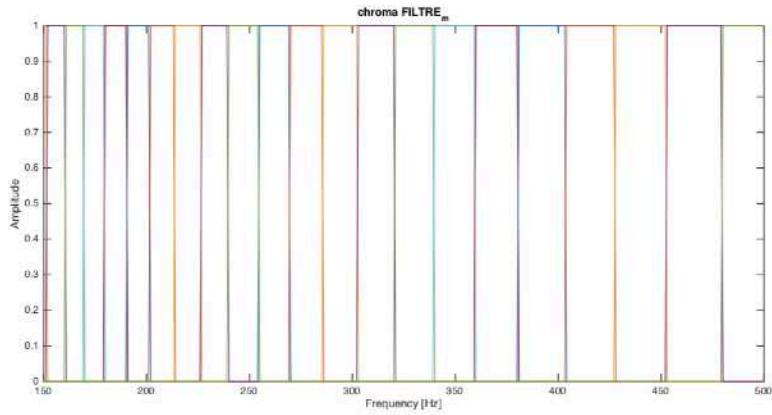
- **Chroma computation**  $C(c, n)$

- We sum up the values of the spectrum  $X(k, n)$  for all  $f_k$  which correspond to a given  $c$
- 1) Hard-mapping
- 2) Soft-mapping

# Chroma/ Pitch Class Profile (PCP)

- **1) Hard-mapping ?**

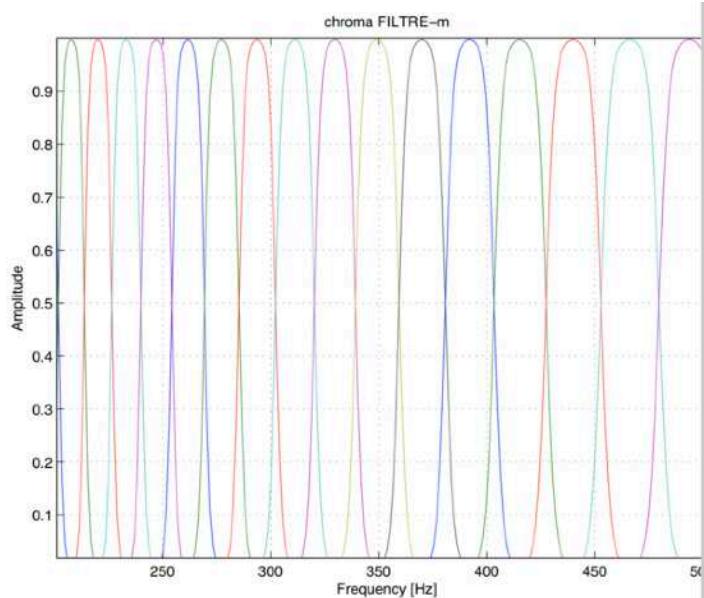
- A frequency  $f_k$  of the DFT only contributes to the closest pitch
- Example:
  - the energy at  $f_k = 452 \text{ Hz}$  ( $p(f_k) = 69.4658$ ) only contributes to the pitch  $p=69$  ( $c=10$ )
  - while  $f_k = 453 \text{ Hz}$  ( $p(f_k) = 69.5041$ ) to  $p=70$  ( $c=11$ ).
- Creation of bank of filters  $H_{p'}$  centered on the semi-tone pitches  $p' \in \{43, 44, \dots, 95\}$ :



# Chroma/ Pitch Class Profile (PCP)

## • 2) Soft-mapping ?

- A frequency  $f_k$  of the DFT contributes to different chromas with a weight inversely proportional to the distance between  $p(f_k)$  and the  $p$  the closest
- Example:
  - the energy at  $f_k = 452 \text{ Hz}$  ( $p(f_k) = 69.4658$ ) contributes nearly equally to  $p=69$  ( $c=10$ ) and  $p=70$  ( $c=11$ ).
- Creation of bank of filters  $H_{p'}$  centered on the semi-tone pitches  $p' \in \{43, 44, \dots, 95\}$ :
  - Each filter is defined by the function
    - $$H_{p'}(f_k) = \frac{1}{2} \tanh(\pi(1 - 2x)) + \frac{1}{2}$$
  - where
    - $x$  = relative distance between the center of the filter  $p'$  and the frequencies of the DFT  $p(f_k)$
    - $$x = R |p' - p(f_k)|$$
  - The filters are evenly distributed and symmetrical on the logarithmic scale of semi-tone pitches, non-zero between  $p' - 1$  and  $p' + 1$  with a maximum value at  $p'$



# Chroma/ Pitch Class Profile (PCP)

- **2) Soft-mapping (cont.)**

- The value of the semi-tone pitch spectrum  $N(n')$  is given by multiplying the values of the DFT  $A(f_k)$  with the bank of filters  $H_{n'}$ :

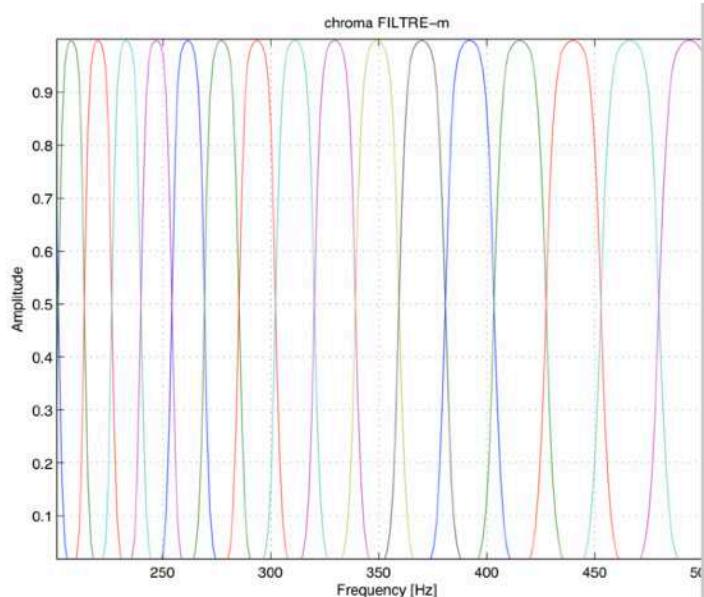
$$P(p') = \sum_{f_k} H_{p'}(f_k)A(f_k)$$

- The mapping between semi-tone pitches  $n$  and the pitch-classes (chroma)  $c$  is defined by:

- $c(p) = \text{mod}(p, 12)$

- The value of the chroma is obtained by summing up the values of equivalent semi-tone pitches

- $C(c) = \sum_{p' \text{ tel que } c(p')=l} P(n') \quad c \in [0, 12[$



# Chroma/ Pitch Class Profile (PCP)

- **Limitations of Chromas - Pitch Class Profile (PCP)**

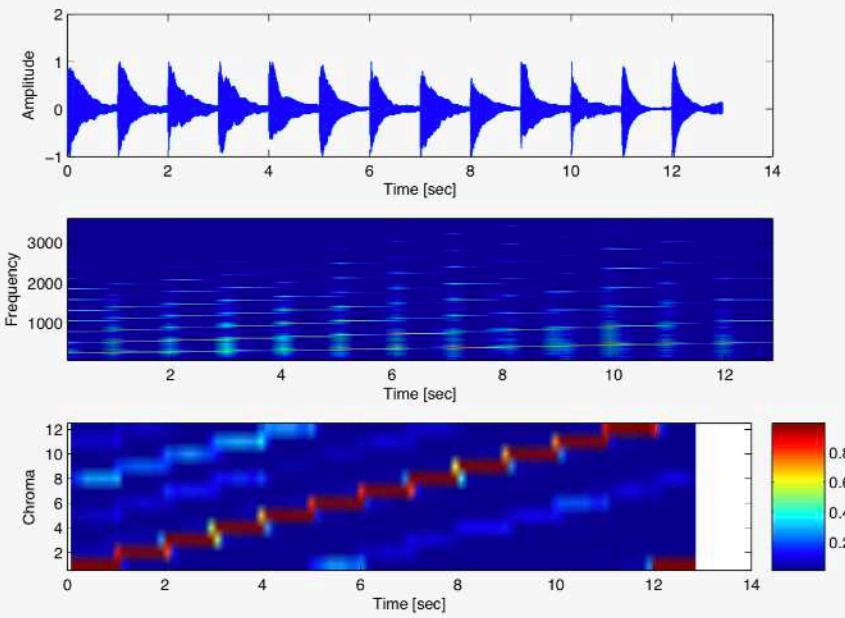
- Presence of the upper harmonics of each note
  - In practice, for a given note  $C$  we don't have  $[1,0,0,0,0,0,0,0,0,0]$
  - but rather  $[a_1 + a_2 + a_4, 0, 0, 0, a_5, 0, 0, a_4, 0, 0, 0, 0]$
  - Influence of the spectral envelope

Pitch	Harmonic	Frequency $f_\mu$	MIDI-scale $m_\mu$	Chroma/PCP $p$
c3	$f_0$	130.81	48	1 (=c)
	$2f_0$	261.62	60	1 (=c)
	$3f_0$	392.43	67.01	8.01 ( $\simeq g$ )
	$4f_0$	523.25	72	1 (=c)
	$5f_0$	654.06	75.86	4.86 ( $\simeq e$ )
	...	...	...	...

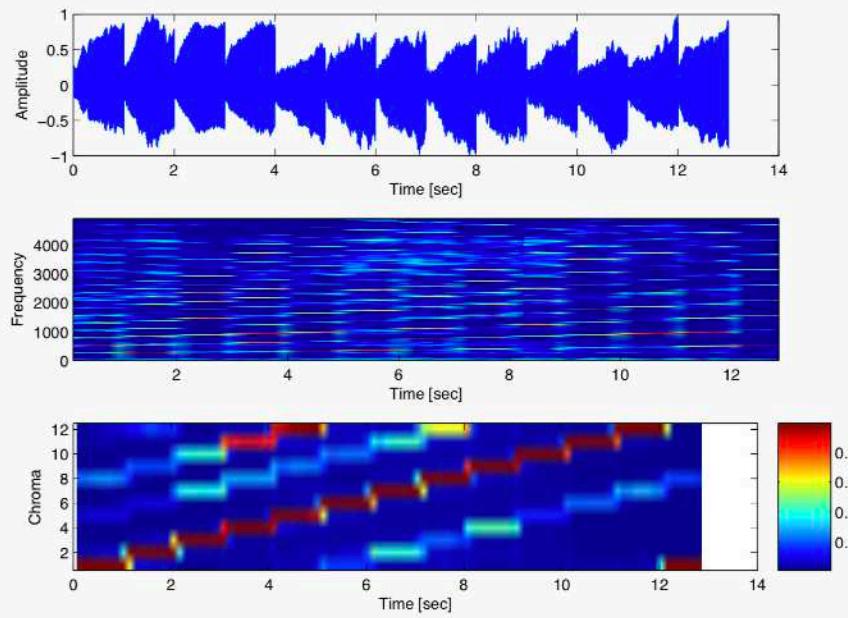
# Chroma/ Pitch Class Profile (PCP)

- **Limitations of Chromas - Pitch Class Profile (PCP)**

Exemple piano

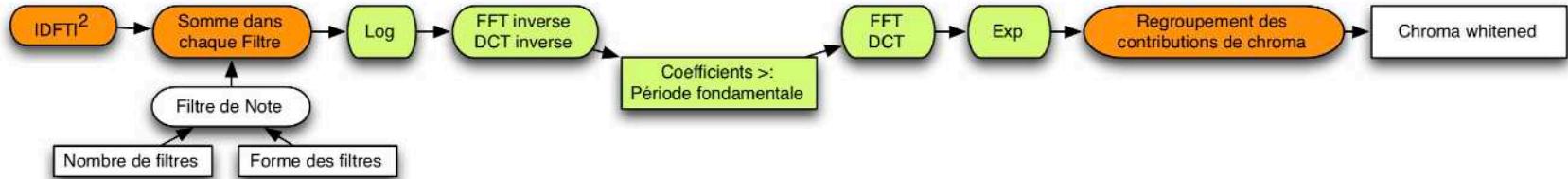
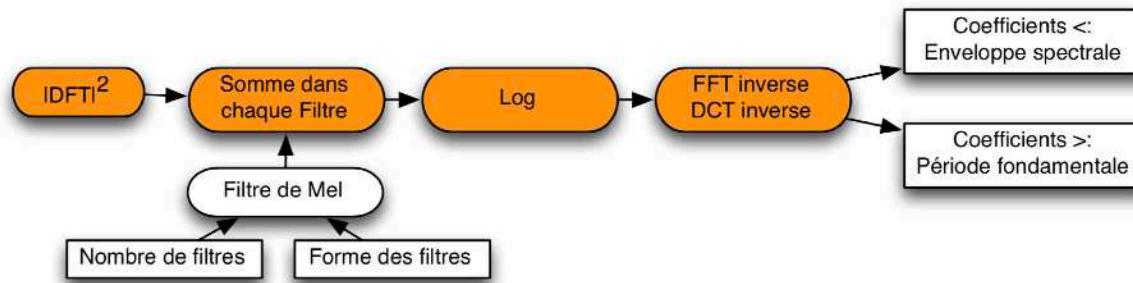
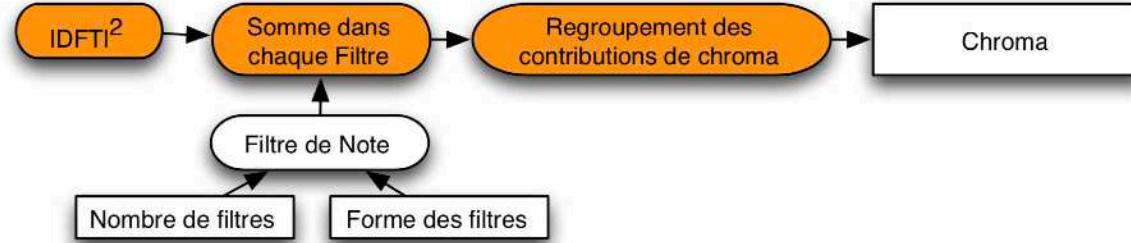


Exemple violon

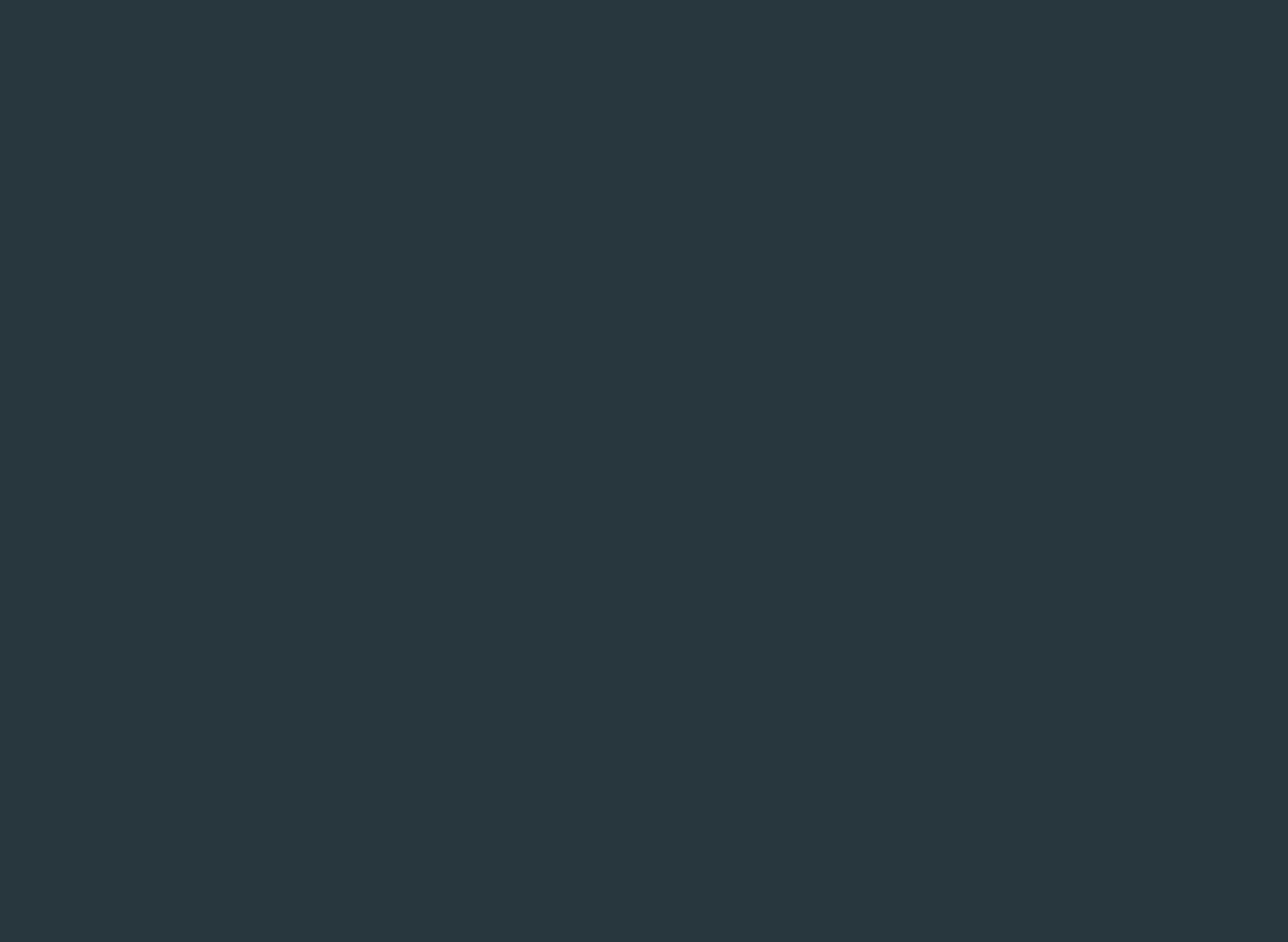


# Chroma/ Pitch Class Profile (PCP)

- Variation of the computation: whitening



# Constant-Q-Transform (CQT)



# Constant-Q-Transform (CQT)

- Discrete Fourier Transform (DFT)

- \_ Definition : Spectral **precision** :  $\Delta f = \frac{sr}{N}$

- it is the step-size at which the Fourier spectrum is sampled
    - it depends on the size of the DFT:  $N$
    - we can improve the precision by increasing  $N$

- \_ Definition : Spectral **resolution** :  $B_w = \frac{C_w}{L}$

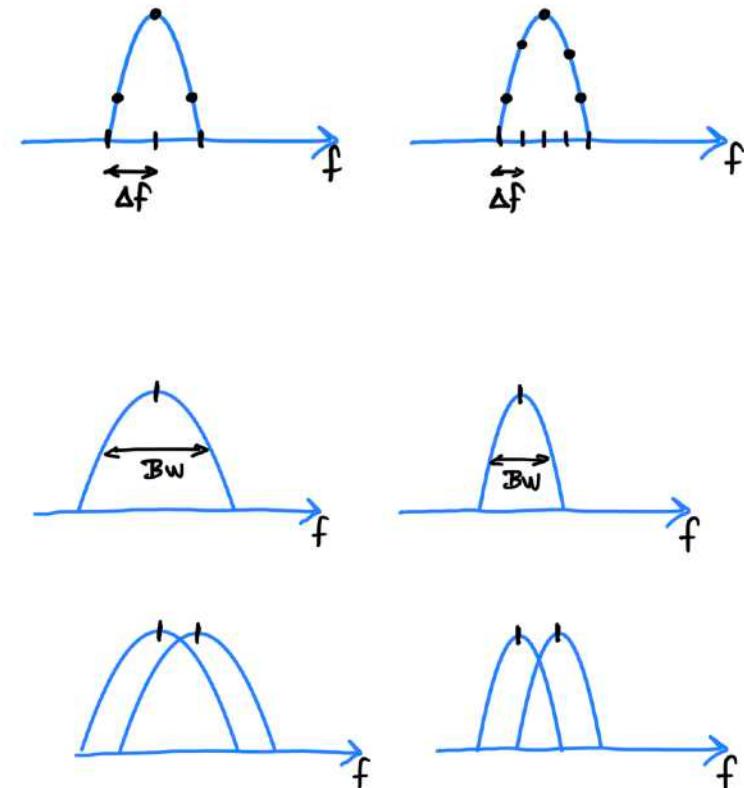
- it describes the ability to discriminate (separate in the spectrum) two adjacent simultaneous frequencies

- Warning :

- even if we increase  $N$  (zero-padding) while keeping  $L$  constant will not improve the resolution !

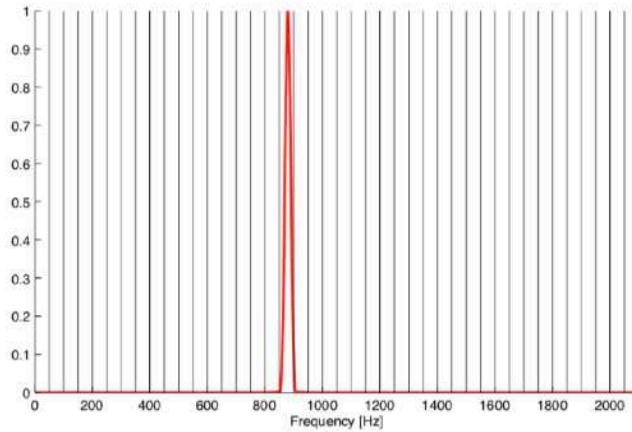
- In a DFT:

- Spectral precision and resolution are constant over frequencies

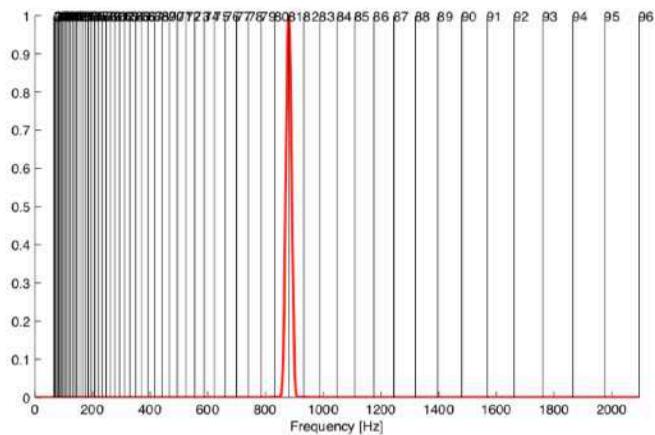


# Constant-Q-Transform (CQT)

- In musical audio
  - the frequencies of the pitches are logarithmically spaced  $f_k = f_0 \cdot 2^{\frac{k}{12}}$ 
    - if we choose A-4 = la-3 = 440 Hz as the reference
    - to go from midi-pitches  $m_k$  to frequencies  $f_k$ :
      - $f_k = 440 \cdot 2^{\frac{m_k - 69}{12}}$
    - to go from frequencies  $f_k$  to midi-pitches  $m_k$ :
      - $m_k = 12 \log_2 \frac{f_k}{440} + 69$
  - pitch frequencies are
    - close together in low frequencies,
    - distant in high frequencies
  - The **spectral resolution** of the DFT
    - is not sufficient (to separate adjacent notes) in low frequencies
    - is too large for high frequencies



Espacement linéaire de la DFT

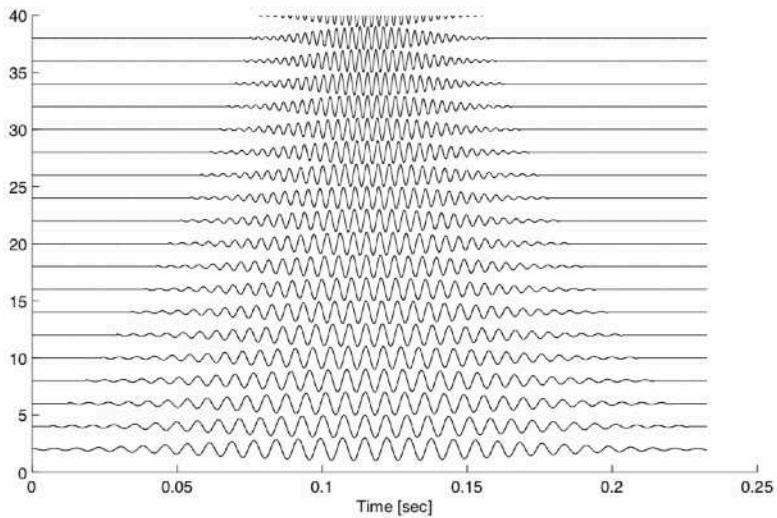


Espacement logarithmique des hauteurs de notes

# Constant-Q-Transform (CQT)

- Solution ?
  - Change the spectral resolution  $B_w$  depending on the frequency  $f_k$  being considered
- How ?
  - By changing the window length  $L$  for each frequency  $f_k$ 
    - The factor  $Q = \frac{f_k}{f_{k+1} - f_k}$  should remains constant in frequency
      - $Q = \frac{f_k}{Bw} = \frac{f_k}{Cw/L} = \frac{f_k \cdot L}{Cw}$
    - We choose a different  $L$  for each frequency  $f_k$ 
      - $L_k = \frac{Q \cdot Cw}{f_k}$

[J. Brown and M. Puckette. An efficient algorithm for the calculation of a constant q transform. JASA, 1992.]



## An efficient algorithm for the calculation of a constant Q transform

John C. Brown  
Hyper-Dimensional Multidimensional Biology, Measurement, and Media Laboratory

Massachusetts Institute of Technology Cambridge Massachusetts 02139

Mark S. Puckette  
CIRCUIT Project MIT Media Lab

(Received 1 February 1992; revised 1 May 1992; accepted 10 June 1992)

**Abstract.** A new algorithm for calculating the discrete Fourier transform (DFT) into a constant  $Q$  transform, where  $Q$  is the ratio of center frequency to bandwidth, has been devised. This algorithm reduces the calculation of terms that are non-adjacent in the corresponding DFT, and it also reduces the number of terms that are non-adjacent in the corresponding  $Q$  transform, while it also adds a small increase to the computation. In effect, this method makes it easier to calculate the spectrum of a signal using a constant  $Q$  transform, and it also makes it easier to calculate the spectrum of a signal using a fast Fourier transform (FFT). Graphical examples of the application of the calculation to musical signals are given for musical pieces by a classical and a rock.

PACS numbers: 43.60.Dv, 43.75.Pq, 43.75.Dm, 43.75.Ez

### 1. INTRODUCTION

In many cases, such as that of musical signals, a constant  $Q$  transform gives a better representation of spectral data than a more commonly employed linear Fourier transform. A recent extension of the constant  $Q$  transform to the nonstationary case has been developed by G. Bäckström (1989), "Semi-constant  $Q$  Transform," and by J. Brown and M. Puckette (1992), "An efficient algorithm for the calculation of a constant  $Q$  transform." Both of these papers have been cited in the literature, and the reader is referred to them for a detailed discussion of the constant  $Q$  transform.

We have calculated a constant  $Q$  transform based on the constant  $Q$  transform of G. Bäckström (1989) in this document. The FFT is calculated using a recursive FFT algorithm, and the result calculated using the slightly longer window lengths. The results are very similar to those obtained by G. Bäckström (1989) for the constant  $Q$  transform.

We have calculated a constant  $Q$  transform based on the constant  $Q$  transform of J. Brown and M. Puckette (1992) in this document. The results are very similar to those obtained by J. Brown and M. Puckette (1992) for the constant  $Q$  transform.

The reason for this is that Eq. (1) is computationally inefficient. However, we can show that the two are identical.

$$\sum_{n=0}^{N-1} x_n e^{-j2\pi f_n n} = \sum_{n=0}^{N-1} x_n e^{-j2\pi f_0 n} e^{-j2\pi f_0 (f_n - f_0)n}, \quad (1)$$

where  $x_n$  is the  $n$ th component of the signal  $x$ ,  $f_n$  is the  $n$ th component of the frequency spectrum, and  $f_0$  is the center frequency.  $x_n e^{-j2\pi f_0 n}$  is a window function of length  $N$ . The expression is the effect of a filter for center frequency  $f_0$ .

In constant  $Q$  transform the center frequencies are proportional to the frequency spectrum. The constant  $Q$  transform is often based on the representation of the signal spectrum as

$$x_n = \sum_{k=0}^{N-1} c_k e^{-j2\pi f_k n}, \quad (2)$$

for some mapping.  $c_k$  is defined as  $c_k = 1/\sqrt{Q}$ , where  $Q$  denotes bandwidth and  $f_k$  denotes frequency. In the case of the filter in Eq. (1), this bandwidth depends on the filter

length  $N$  (Brown, 1992) for resolution.

We can use Eq. (2) to evaluate Eq. (1) as follows. Let

$$\langle x_n | x_m \rangle = \sum_{k=0}^{N-1} c_k \langle x_k | x_m \rangle e^{-j2\pi f_k m}, \quad (3)$$

where  $\langle x_n | x_m \rangle$  is the  $n$ th component of the constant  $Q$  transform of  $x$ , and  $\langle x_k | x_m \rangle$  is the  $k$ th component of the constant  $Q$  transform of  $x$ . Equation (3) is a form of Parseval's equation (Oppenheim et al., 1987).

The term  $\langle x_n | x_m \rangle$  is the discrete Fourier transform of  $x$  at frequency  $f_m$  (Brown, 1992).

Equation (3) gives

$$\sum_{n=0}^{N-1} x_n \overline{x_m} = \sum_{k=0}^{N-1} c_k \langle x_k | x_m \rangle, \quad (4)$$

where  $\overline{x_m}$  is the discrete Fourier transform of  $x$  at frequency  $f_m$  (Brown, 1992).

$$\sum_{n=0}^{N-1} x_n \overline{x_m} = \sum_{k=0}^{N-1} c_k \sum_{n=0}^{N-1} x_n e^{-j2\pi f_k n} e^{j2\pi f_m n}, \quad (5)$$

where  $\langle x_k | x_m \rangle$  is the discrete Fourier transform of  $x$  at frequency  $f_m$  (Brown, 1992).

$$\sum_{n=0}^{N-1} x_n \overline{x_m} = \sum_{k=0}^{N-1} c_k \sum_{n=0}^{N-1} x_n e^{-j2\pi f_k n} e^{-j2\pi f_0 (f_m - f_0)n}, \quad (6)$$

We will note that  $\langle x_k | x_m \rangle$  is the frequency domain window function of length  $N$  (Brown, 1992).

We have a filtering effect on the signal spectrum due to the window function.

**2. THEORETICAL**

The theoretical derivation of the constant  $Q$  transform is based on the following assumptions:

(1) The signal is a stationary random process.

(2) The signal is a narrow-band signal.

(3) The signal is a complex exponential signal.

(4) The signal is a pure tone.

(5) The signal is a periodic signal.

(6) The signal is a periodic signal.

(7) The signal is a periodic signal.

(8) The signal is a periodic signal.

(9) The signal is a periodic signal.

(10) The signal is a periodic signal.

(11) The signal is a periodic signal.

(12) The signal is a periodic signal.

(13) The signal is a periodic signal.

(14) The signal is a periodic signal.

(15) The signal is a periodic signal.

(16) The signal is a periodic signal.

(17) The signal is a periodic signal.

(18) The signal is a periodic signal.

(19) The signal is a periodic signal.

(20) The signal is a periodic signal.

(21) The signal is a periodic signal.

(22) The signal is a periodic signal.

(23) The signal is a periodic signal.

(24) The signal is a periodic signal.

(25) The signal is a periodic signal.

(26) The signal is a periodic signal.

(27) The signal is a periodic signal.

(28) The signal is a periodic signal.

(29) The signal is a periodic signal.

(30) The signal is a periodic signal.

(31) The signal is a periodic signal.

(32) The signal is a periodic signal.

(33) The signal is a periodic signal.

(34) The signal is a periodic signal.

(35) The signal is a periodic signal.

(36) The signal is a periodic signal.

(37) The signal is a periodic signal.

(38) The signal is a periodic signal.

(39) The signal is a periodic signal.

(40) The signal is a periodic signal.

(41) The signal is a periodic signal.

(42) The signal is a periodic signal.

(43) The signal is a periodic signal.

(44) The signal is a periodic signal.

(45) The signal is a periodic signal.

(46) The signal is a periodic signal.

(47) The signal is a periodic signal.

(48) The signal is a periodic signal.

(49) The signal is a periodic signal.

(50) The signal is a periodic signal.

(51) The signal is a periodic signal.

(52) The signal is a periodic signal.

(53) The signal is a periodic signal.

(54) The signal is a periodic signal.

(55) The signal is a periodic signal.

(56) The signal is a periodic signal.

(57) The signal is a periodic signal.

(58) The signal is a periodic signal.

(59) The signal is a periodic signal.

(60) The signal is a periodic signal.

(61) The signal is a periodic signal.

(62) The signal is a periodic signal.

(63) The signal is a periodic signal.

(64) The signal is a periodic signal.

(65) The signal is a periodic signal.

(66) The signal is a periodic signal.

(67) The signal is a periodic signal.

(68) The signal is a periodic signal.

(69) The signal is a periodic signal.

(70) The signal is a periodic signal.

(71) The signal is a periodic signal.

(72) The signal is a periodic signal.

(73) The signal is a periodic signal.

(74) The signal is a periodic signal.

(75) The signal is a periodic signal.

(76) The signal is a periodic signal.

(77) The signal is a periodic signal.

(78) The signal is a periodic signal.

(79) The signal is a periodic signal.

(80) The signal is a periodic signal.

(81) The signal is a periodic signal.

(82) The signal is a periodic signal.

(83) The signal is a periodic signal.

(84) The signal is a periodic signal.

(85) The signal is a periodic signal.

(86) The signal is a periodic signal.

(87) The signal is a periodic signal.

(88) The signal is a periodic signal.

(89) The signal is a periodic signal.

(90) The signal is a periodic signal.

(91) The signal is a periodic signal.

(92) The signal is a periodic signal.

(93) The signal is a periodic signal.

(94) The signal is a periodic signal.

(95) The signal is a periodic signal.

(96) The signal is a periodic signal.

(97) The signal is a periodic signal.

(98) The signal is a periodic signal.

(99) The signal is a periodic signal.

(100) The signal is a periodic signal.

(101) The signal is a periodic signal.

(102) The signal is a periodic signal.

(103) The signal is a periodic signal.

(104) The signal is a periodic signal.

(105) The signal is a periodic signal.

(106) The signal is a periodic signal.

(107) The signal is a periodic signal.

(108) The signal is a periodic signal.

(109) The signal is a periodic signal.

(110) The signal is a periodic signal.

(111) The signal is a periodic signal.

(112) The signal is a periodic signal.

(113) The signal is a periodic signal.

(114) The signal is a periodic signal.

(115) The signal is a periodic signal.

(116) The signal is a periodic signal.

(117) The signal is a periodic signal.

(118) The signal is a periodic signal.

(119) The signal is a periodic signal.

(120) The signal is a periodic signal.

(121) The signal is a periodic signal.

(122) The signal is a periodic signal.

(123) The signal is a periodic signal.

(124) The signal is a periodic signal.

(125) The signal is a periodic signal.

(126) The signal is a periodic signal.

(127) The signal is a periodic signal.

(128) The signal is a periodic signal.

(129) The signal is a periodic signal.

(130) The signal is a periodic signal.

(131) The signal is a periodic signal.

(132) The signal is a periodic signal.

(133) The signal is a periodic signal.

(134) The signal is a periodic signal.

(135) The signal is a periodic signal.

(136) The signal is a periodic signal.

(137) The signal is a periodic signal.

(138) The signal is a periodic signal.

(139) The signal is a periodic signal.

(140) The signal is a periodic signal.

(141) The signal is a periodic signal.

(142) The signal is a periodic signal.

(143) The signal is a periodic signal.

(144) The signal is a periodic signal.

(145) The signal is a periodic signal.

(146) The signal is a periodic signal.

(147) The signal is a periodic signal.

(148) The signal is a periodic signal.

(149) The signal is a periodic signal.

(150) The signal is a periodic signal.

(151) The signal is a periodic signal.

(152) The signal is a periodic signal.

(153) The signal is a periodic signal.

(154) The signal is a periodic signal.

(155) The signal is a periodic signal.

(156) The signal is a periodic signal.

(157) The signal is a periodic signal.

(158) The signal is a periodic signal.

(159) The signal is a periodic signal.

(160) The signal is a periodic signal.

(161) The signal is a periodic signal.

(162) The signal is a periodic signal.

(163) The signal is a periodic signal.

(164) The signal is a periodic signal.

(165) The signal is a periodic signal.

(166) The signal is a periodic signal.

(167) The signal is a periodic signal.

(168) The signal is a periodic signal.

(169) The signal is a periodic signal.

(170) The signal is a periodic signal.

(171) The signal is a periodic signal.

(172) The signal is a periodic signal.

(173) The signal is a periodic signal.

(174) The signal is a periodic signal.

(175) The signal is a periodic signal.

(176) The signal is a periodic signal.

(177) The signal is a periodic signal.

(178) The signal is a periodic signal.

(179) The signal is a periodic signal.

(180) The signal is a periodic signal.

(181) The signal is a periodic signal.

(182) The signal is a periodic signal.

(183) The signal is a periodic signal.

(184) The signal is a periodic signal.

(185) The signal is a periodic signal.

(186) The signal is a periodic signal.

(187) The signal is a periodic signal.

(188) The signal is a periodic signal.

(189) The signal is a periodic signal.

(190) The signal is a periodic signal.

(191) The signal is a periodic signal.

(192) The signal is a periodic signal.

(193) The signal is a periodic signal.

(194) The signal is a periodic signal.

(195) The signal is a periodic signal.

(196) The signal is a periodic signal.

(197) The signal is a periodic signal.

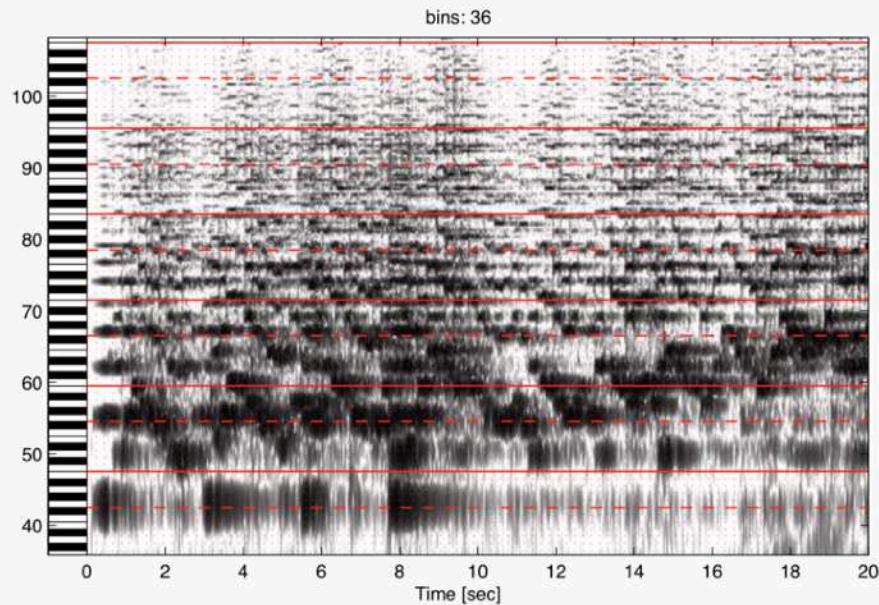
(198) The signal is a periodic signal.

(199) The signal is a periodic signal.

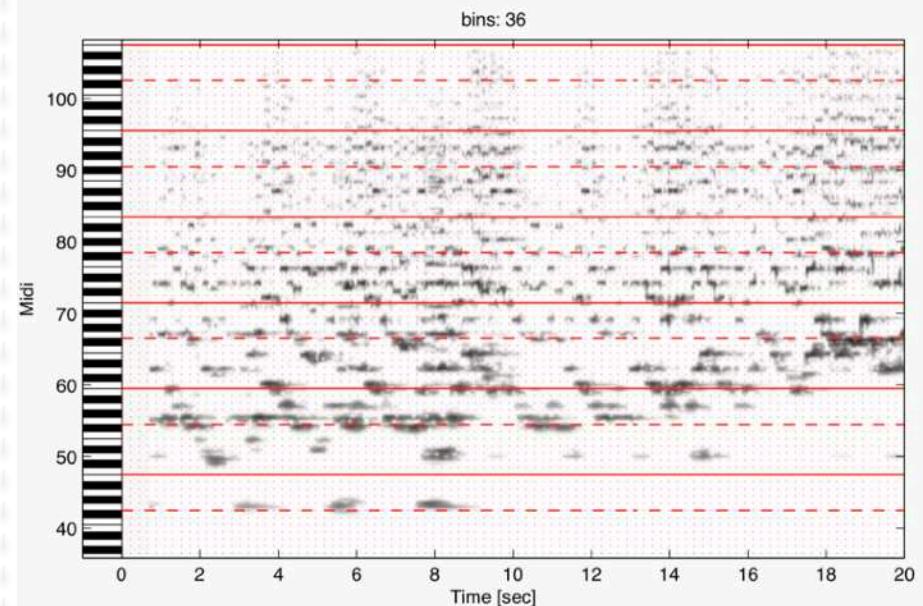
(200) The signal is a periodic signal.

# Constant-Q-Transform (CQT)

Example (using DFT)

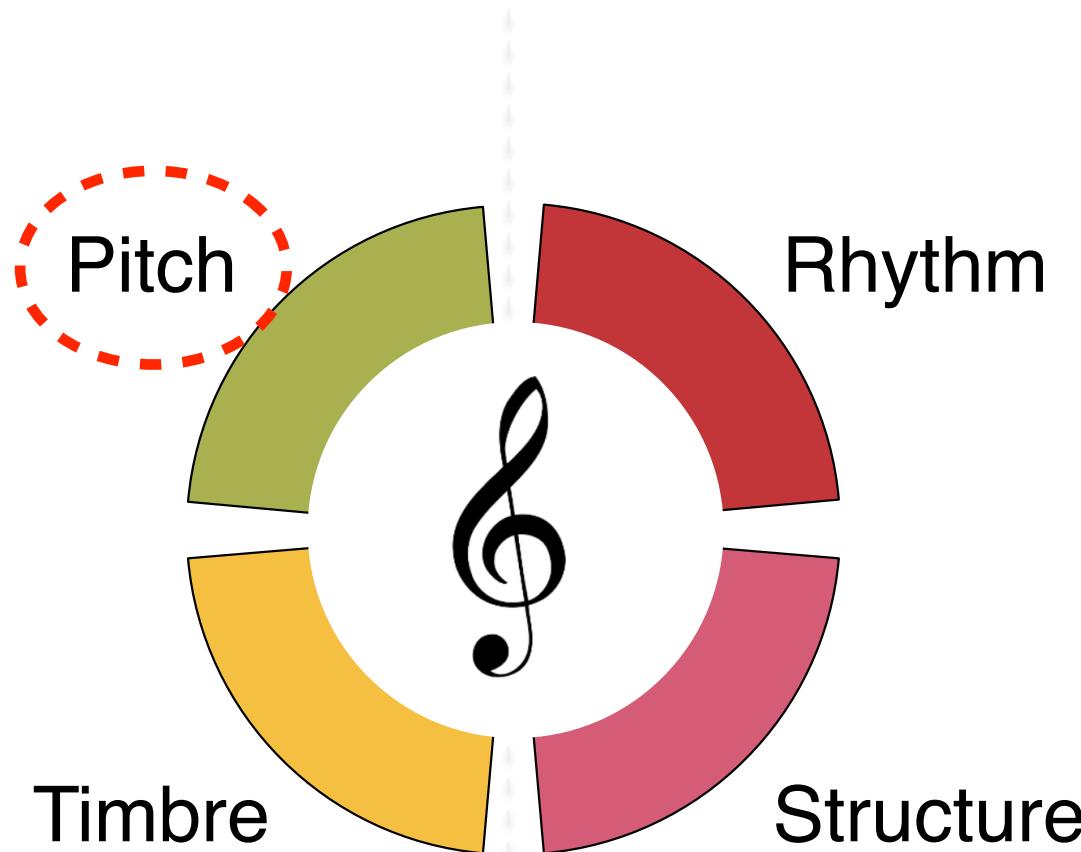


Example (using the CQT)



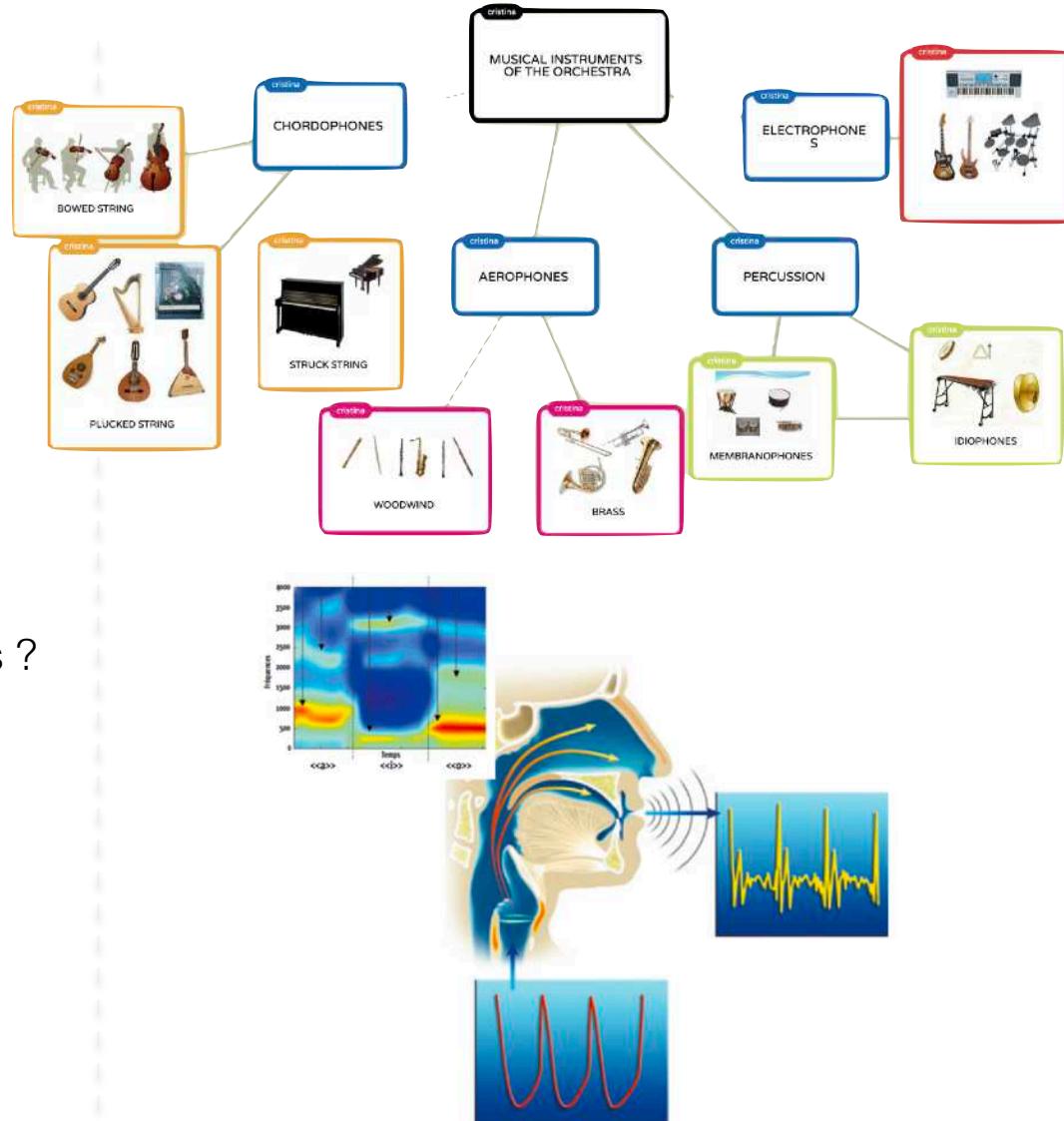
# Reminder of Musical Concepts

# What is pitch ?



# Harmonic sounds

- Most musical instruments produce harmonic sounds
  - strings (plucked, bowed, hammered)
  - woodwinds, brass



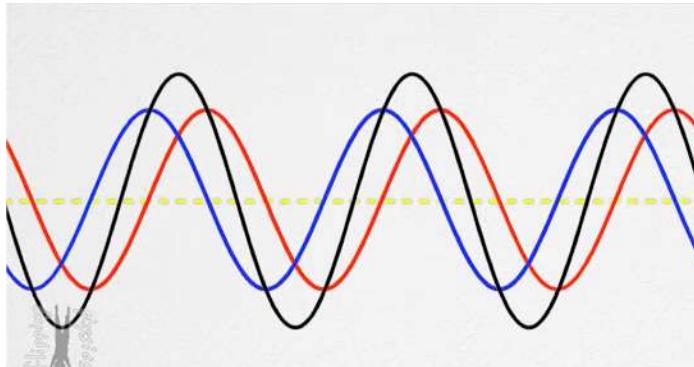
- Why do they produce harmonic sounds ?
  - Two visions:
    - 1) periodic signal
    - 2) vibration modes

# Harmonic sounds

## Standing waves

- **Reminder**

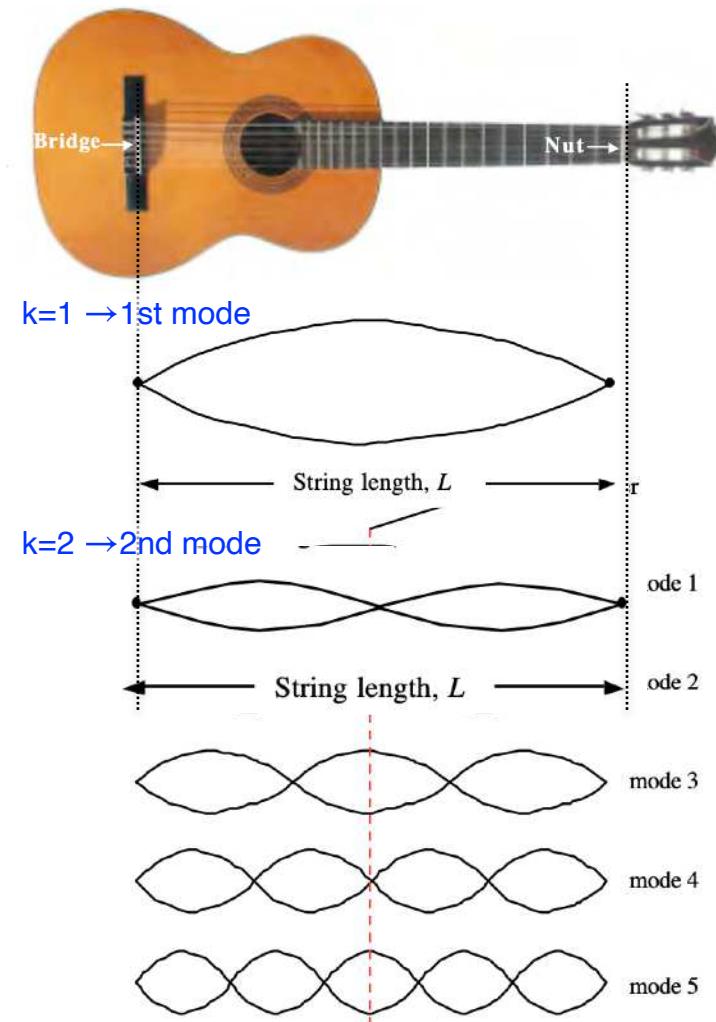
- **Standing wave** in a string fixed at both ends



- Possible **frequency modes** of a fingered string

$$f_k = k \frac{v}{\lambda} = k \frac{v}{2L} = k \frac{\sqrt{T/\mu}}{2L}$$

- $k$  vibration mode or harmonic number
  - $T$ : string tension,
  - $\mu$ : linear mass density,



# Harmonic sounds

## Fourier series

- Harmonic representation of sound
  - Fourier series

$$f(t) = \sum_{n=0}^{+\infty} k_n \sin(n\omega_0 t + \phi_n)$$

$\sin(a + b) = \sin(a)\cos(b) + \cos(a)\sin(b)$

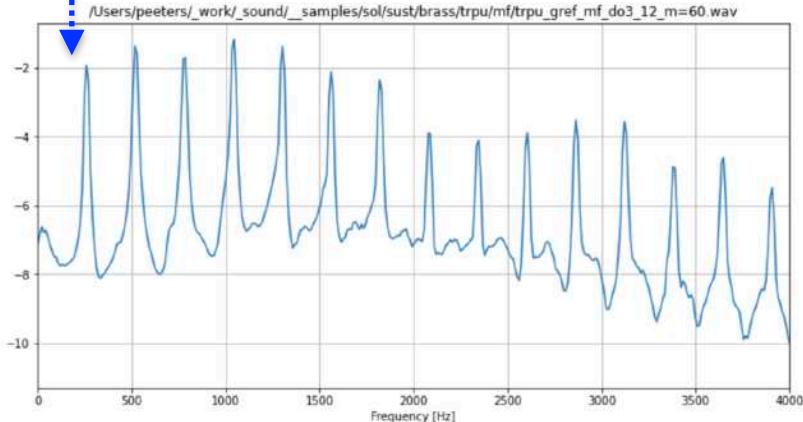
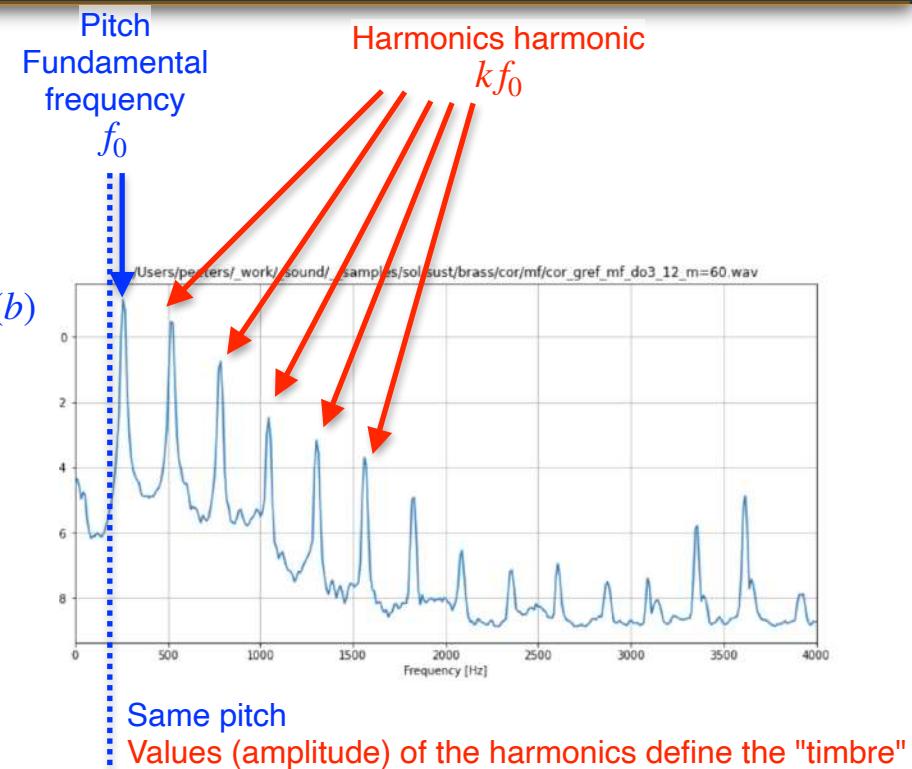
$$= \sum_{n=0}^{+\infty} A_n \cos(n\omega_0 t) + B_n \sin(n\omega_0 t)$$

$$= \sum_{n=-\infty}^{+\infty} c_n e^{jn\omega_0 t}$$

$$c_n = \begin{cases} \frac{A_n + iB_n}{2} & \text{if } n < 0 \\ \frac{A_0}{2} & \text{if } n = 0 \\ \frac{A_n - iB_n}{2} & \text{if } n > 0 \end{cases}$$

- Pitch (fundamental frequency)

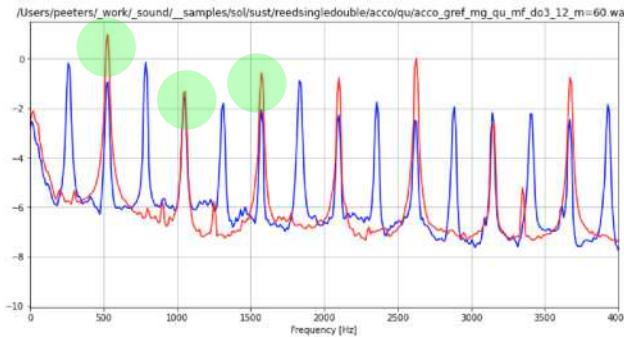
Mode of vibration	Frequency name (for any type of overtone)	Frequency name (for harmonic overtones)
First	Fundamental	First harmonic
Second	First overtone	Second harmonic
Third	Second overtone	Third harmonic
Fourth	Third overtone	Fourth harmonic



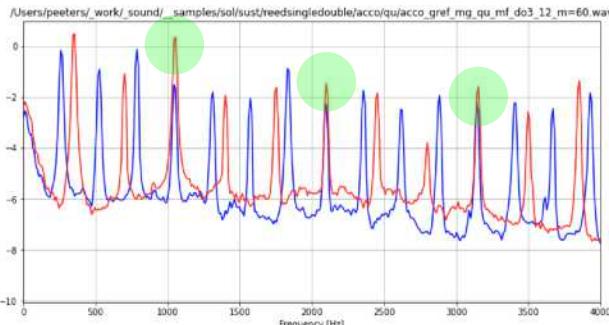
# Harmonic sounds Intervals

- Consonant/ Dissonant intervals between pitches
  - can be explained by the overlap between the harmonics
  - closeness between harmonics that will fall into the same critical band

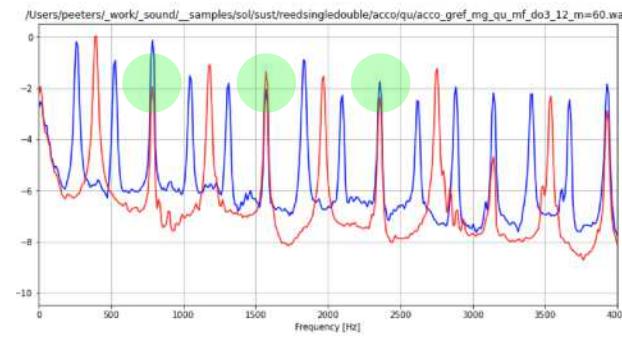
Octave interval ( $F_0 = 2f_0$ )



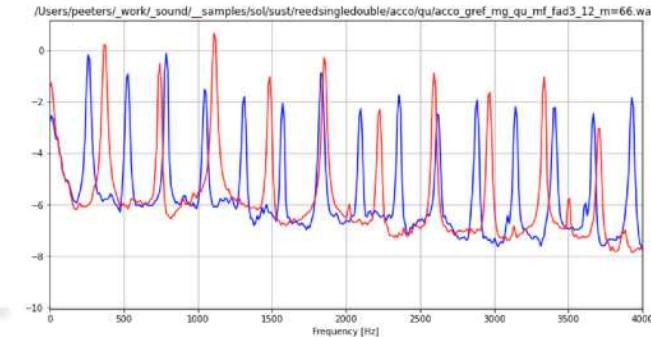
4<sup>th</sup> interval ( $3F_0 = 4f_0$ )



5<sup>th</sup> interval ( $2F_0 = 3f_0$ )



Triton (most dissonant !!!)

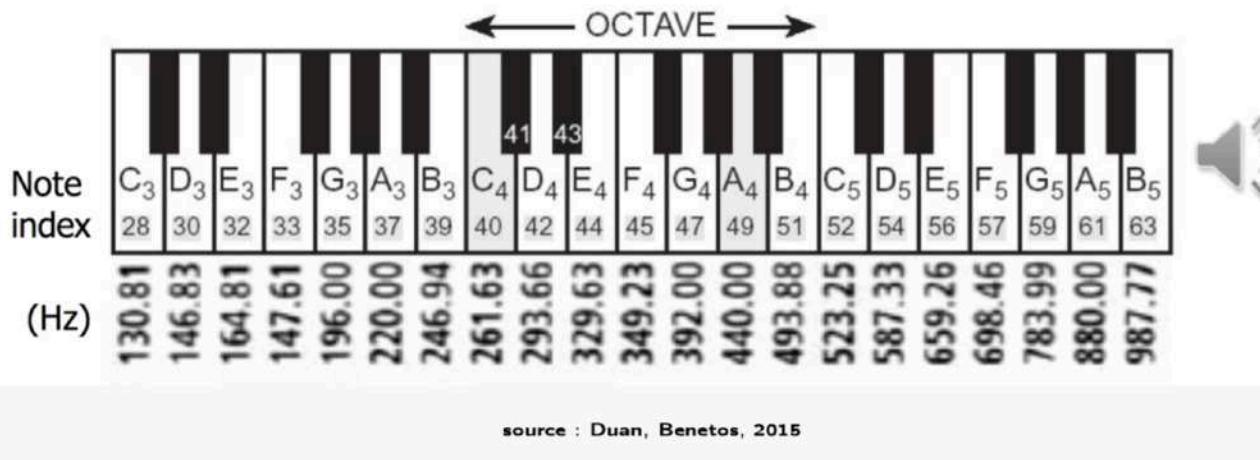


# What is pitch ?

- **Pitch:**

- the attribute of auditory sensation in terms of which sounds may be **ordered** on a scale extending from low to high (ANSI)
- (Operational) A sound has a certain pitch if it can be reliably matched to a **sine tone** of a given frequency at 40 dB SPL
- People hear pitch in a **logarithmic** scale

- **Notes** = Musical representation of pitch: do-3, re-3, mi-3



# What is pitch ?

## Mapping Pitch to Notes

- To **map Pitch to Notes** we need to define an **organisation of the pitches**
  - **Temperament = tuning system** that allows representing (at best) the intervals of just intonation
  - Temperament refers to the various tuning systems for the subdivision of the octave
- **Intervals** ? the difference in pitch between two sounds
  - **Octave:** do-3 → do-4 (most consonant) ... because of the harmonic series ⇒ 2<sup>nd</sup> harmonics
  - **Fifth:** do-3 → sol-3 (very consonant) ... because of the harmonic series ⇒ 3<sup>rd</sup> harmonics
  - **Fourth:** do-3 → fa-3 (very consonant) ... because of the harmonic series ⇒ 4<sup>th</sup> harmonics
- The various **temperament** differs according to **how an octave is subdivided into notes** to allow representing the intervals

# What is pitch ?

## Temperaments

- **(1) Pythagorean temperament**

- Pythagorean temperament = based on natural proportions **3:2**
- Consider a string of length  $L$ , it produces a frequency  $f = k \frac{v}{2L}$
- Divide/Multiply the length by a factor 3:2
- If we finger the string a  $L' =$     **each time we bring it back to main octave  $\in [1,2]$**

- $L \left( \frac{1}{2} \right) \rightarrow 2f$

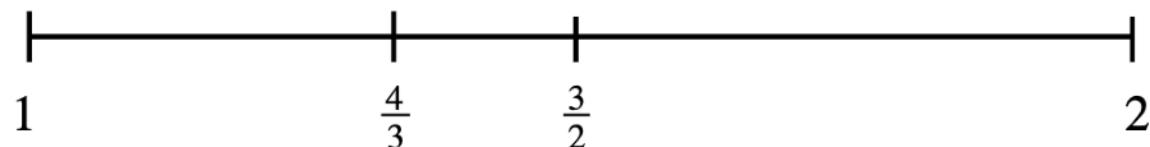
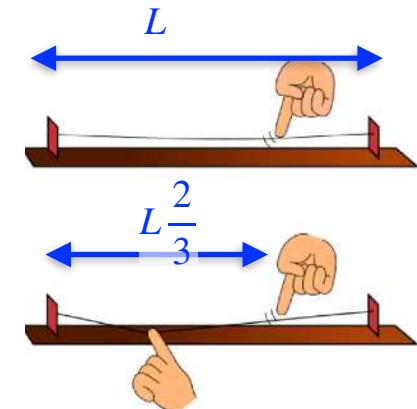
→ octave (consonant interval) = **do**

- $L \left( \frac{2}{3} \right) \rightarrow \frac{3}{2}f$

→ 5th (consonant interval) = **sol**

- $L \left( \frac{3}{2} \right) \rightarrow \frac{2}{3}f \rightarrow \frac{4}{3}f$

→ 4th (consonant interval) = **fa**



# What is pitch ?

## Temperaments

- **(1) Pythagorean temperament**

- If we finger the string a  $L' =$  each time we bring it back to main octave  $\in [1,2]$

- $L \left( \frac{2}{3} \right) \rightarrow \frac{3}{2} f = \text{sol}$

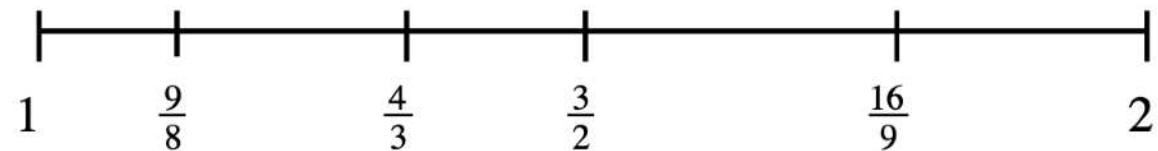
$$L \left( \frac{2}{3} \right)^2 \rightarrow \frac{3^2}{2^2} f \rightarrow \frac{9}{8} f = \text{ré}$$

- $L \left( \frac{3}{2} \right) \rightarrow \frac{2}{3} f \rightarrow \frac{4}{3} f = \text{fa}$

$$L \left( \frac{3}{2} \right)^2 \rightarrow \frac{2^2}{3^2} f \rightarrow \frac{16}{9} f = \text{si-b}$$

-

Pentatonic scale !



# What is pitch ?

## Temperaments

- **(1) Pythagorean temperament**

- If we finger the string a  $L' =$

- $L\left(\frac{2}{3}\right) \rightarrow \frac{3}{2}f = \text{sol}$

- $L\left(\frac{3}{2}\right) \rightarrow \frac{2}{3}f \rightarrow \frac{4}{3}f = \text{fa}$

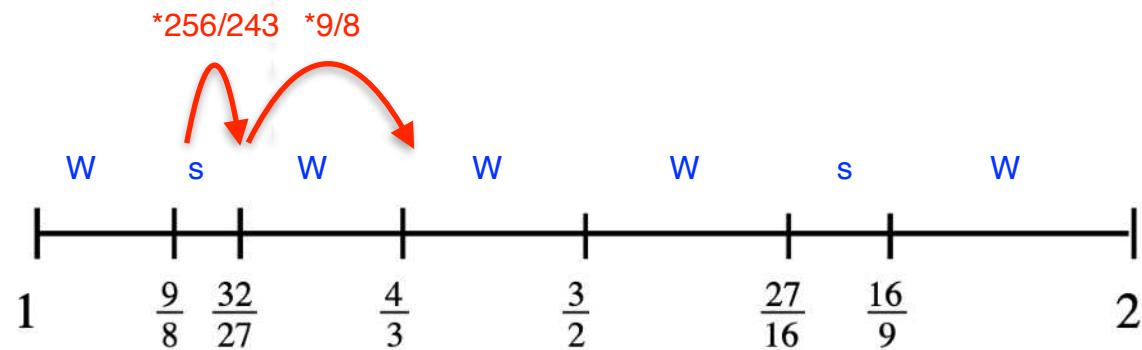
each time we bring it back to main octave  $\in [1,2]$

$$L\left(\frac{2}{3}\right)^2 \rightarrow \frac{3^2}{2^2}f \rightarrow \frac{9}{8}f = \text{ré}$$

$$L\left(\frac{3}{2}\right)^2 \rightarrow \frac{2^2}{3^2}f \rightarrow \frac{16}{9}f = \text{si-b}$$

$$L\left(\frac{2}{3}\right)^3 \rightarrow \frac{3^3}{2^3}f \rightarrow \frac{27}{16}f = \text{la}$$

$$L\left(\frac{3}{2}\right)^3 \rightarrow \frac{2^3}{3^3}f \rightarrow \frac{32}{27}f = \text{mi-b}$$



# What is pitch ?

## Temperaments

- **(1) Pythagorean temperament**

- Pythagorean temperament = based on natural proportions 3:2
- Consider a string of length  $L$ , it produces a frequency  $f = k \frac{v}{2L}$
- If we finger the string a  $L' =$  each time we bring it back to main octave  $\in [1,2]$

$$\bullet L \left( \frac{1}{2} \right) \rightarrow 2f$$

→ octave (consonant interval) = **do**

$$\bullet L \left( \frac{2}{3} \right) \rightarrow \frac{3}{2} f = 1.5f$$

→ 5th (consonant interval) = **sol**

$$\bullet L \left( \frac{2}{3} \right)^2 \rightarrow \frac{3^2}{2^2} f \rightarrow f \frac{3^2}{2^3} f \approx 1.13$$

→ 5th (5th) back to main octave (divide frequency by 2) = **ré**

$$\bullet L \left( \frac{2}{3} \right)^3 \rightarrow \frac{3^3}{2^3} f \rightarrow \frac{3^3}{2^4} f \approx 1.69$$

→ 5th (5th (5th)) back to main octave = **la**

• and so one

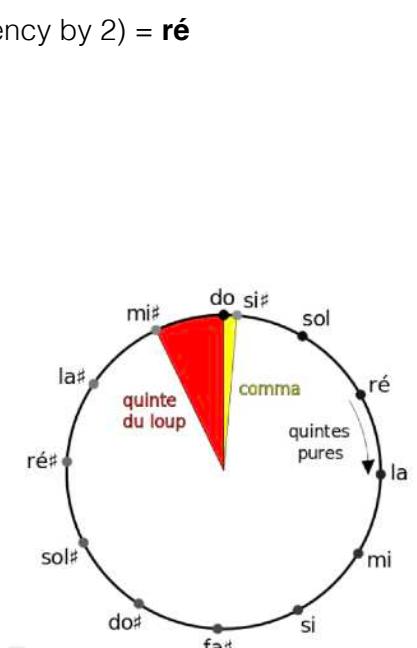
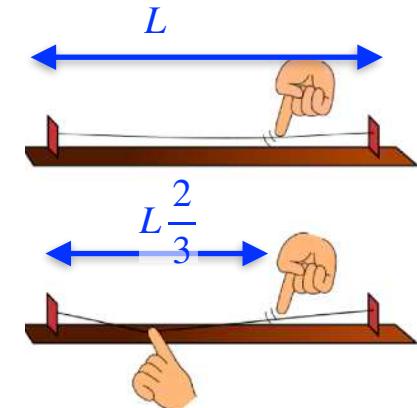
- Resulting scale

- problem: we are not getting back to do:  $\frac{3^{12}}{2^{12}} \frac{1}{2^7} = 1.0136 \neq 1$

- to get back to the starting note (do) we need to use a special fifth

- Intervals are not equal !!!

- not possible to transpose a given piece of music !!!
- semi-tone intervals are not exactly half of the whole tones



# What is pitch ? Temperaments

- **(2) Equal temperament**

- divide an octave in 12 equal intervals →intervals:  $\frac{f_1}{f_0} = 2^{\frac{1}{12}}$

$$f_0 = f \cdot 2^{\frac{0}{12}}$$

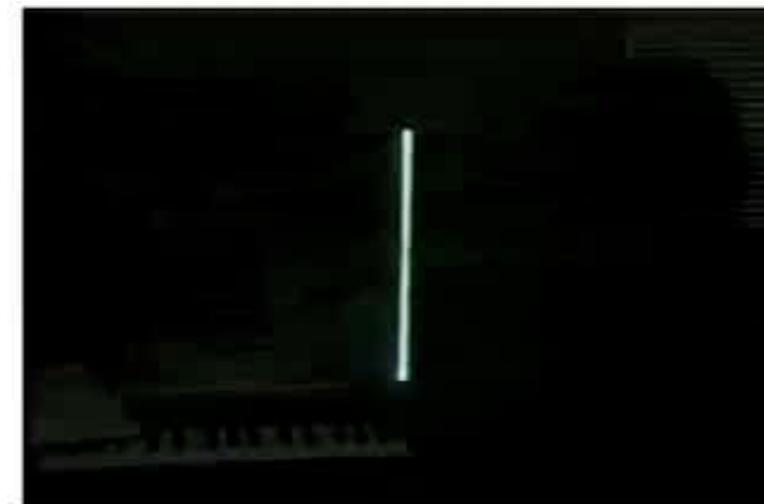
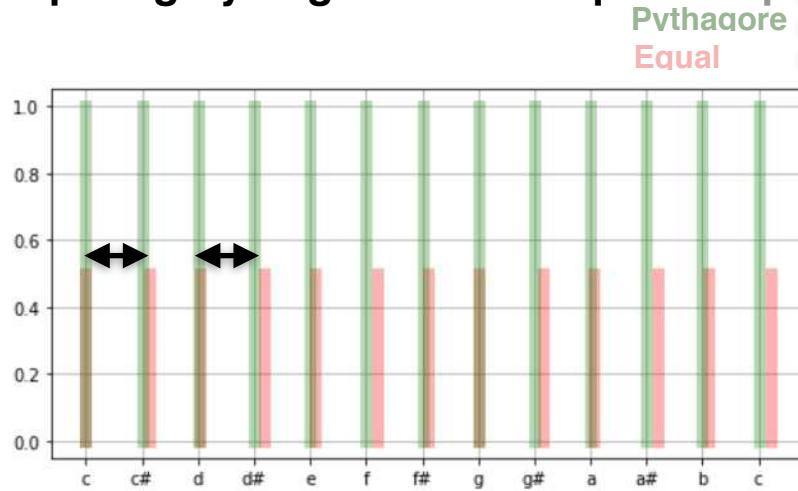
$$f_1 = f \cdot 2^{\frac{1}{12}}$$

$$f_2 = f \cdot 2^{\frac{2}{12}}$$

...

$$f_{12} = f \cdot 2^{\frac{12}{12}} = 2 \cdot f_0$$

- **Comparing Pythagorean and Equal Temperament**



[https://fr.wikipedia.org/wiki/Courbe\\_de\\_Lissajous](https://fr.wikipedia.org/wiki/Courbe_de_Lissajous)

<https://www.youtube.com/watch?v=6NII4No3sOM>

# What is pitch ?

- **Mapping Pitch to Notes** (using the equal temperament)

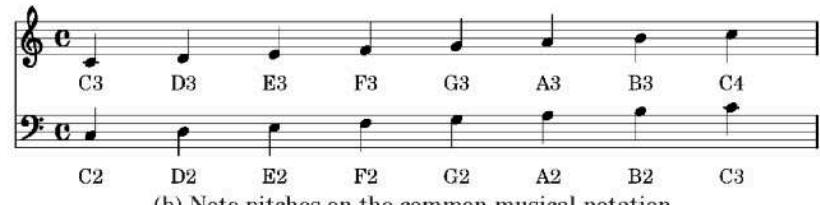
- 261 Hz → C-4 or do-3 (do at 3<sup>rd</sup> octave)
- 440 Hz → A-4 or la-3 (la at 3<sup>rd</sup> octave)
- 880 Hz → A-5 or la-4 (la at 4<sup>th</sup> octave)

- midi values

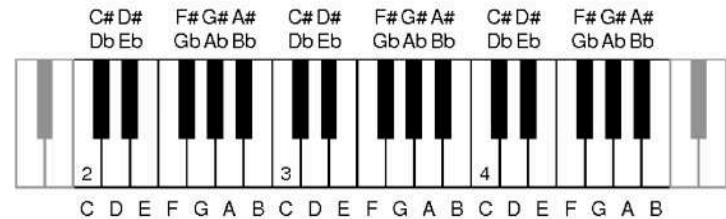
- $f = 440 \cdot 2^{\frac{(m-69)}{12}}$
- $m = 12 \log_2 \left( \frac{f}{440} \right) + 69$
- 261 Hz → 60
- 440 Hz → 69
- 880 Hz → 81

Note $f_0, \text{MIDI}$	C3	D3	E3	F3	G3	A3	B3	C4
$f_0, \text{Hz}$	48	50	52	53	55	57	59	60

(a) Fundamental frequencies and MIDI note numbers for note names.



(b) Note pitches on the common musical notation.



# What is pitch ?

## How pitches are organised ?



<https://www.youtube.com/watch?v=e5Ip742VlyQ>

- **Musical textures**

- **Monophonic**

- one voice, monophonic texture (monophony) refers to a single melodic line, though it may be played by one or many instruments
    - voices may be in exact unison or in different octaves, as long as the same notes and rhythms are played

- **Homophonic**

- one voice, a melody, which stands out from background accompaniment
    - accompaniment may be simple chords or a harmony with melodic interest, but in either case, the main melody must be clearly distinguishable

- **Polyphonic**

- or counterpoint involves multiple melodic voices, all of equal importance, occurring simultaneously.
    - complex, dense texture is typical of Renaissance and baroque music

The diagram illustrates three types of musical textures:

- Monophony:** A single melodic line (red wavy line) over a simple harmonic background (one note per measure).
- Homophony:** A single melodic line (yellow wavy line) supported by a harmonic background (chords).
- Polyphony:** Multiple melodic voices (green and blue wavy lines) occurring simultaneously.

# What is pitch ?

## How pitches are organised ?

- **Modal music**

- since the Greeks, Gregorian and Church
- which uses modal scales
  - (as opposed to Major/minor scales)
- centred on horizontal-melody often monophonic
  - (as opposed to vertical-harmony),
- in other cultures:
  - Indian music ("râgas"), Turkish music ("maqâm"), or modal jazz (Miles Davis "Kind of Blue")

The image shows six musical staves, each with a different mode name above it. The first staff is labeled 'Ionian Major' and has a red box around it. The second staff is 'Dorian Minor #6'. The third staff is 'Phrygian Minor b2'. The fourth staff is 'Lydian Major #4'. The fifth staff is 'Mixolydian Major b7'. The sixth staff is 'Aeolian Minor (natural)' and also has a red box around it.

Greek, Gregorian, Church modes

The image shows nine musical staves, each with a different raga name above it. The first raga is 'Bilâwal'. The second is 'Kalyan'. The third is 'Bhairav'. The fourth is 'Purvi'. The fifth is 'Mârvâ'. The sixth is 'Bhairavi'. The seventh is 'Åsâvâri'. The eighth is 'Kâfi'. The ninth is 'Todî'. Each staff contains a series of notes on a treble clef staff.

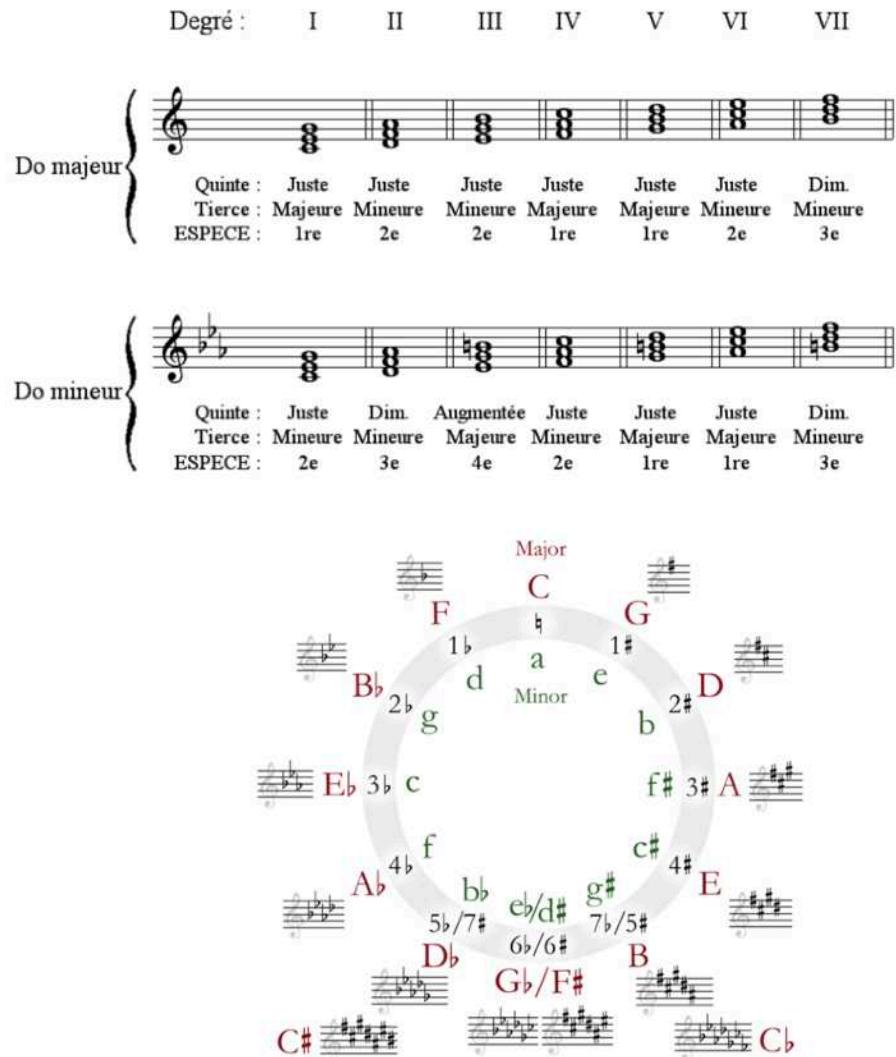
Indian raga

# What is pitch ?

## How pitches are organised ?

- **Tonal music**

- centred around the tonality
- define the set of relationships between notes, structured chords around a given tonality
- tonal language is built upon the diatonic Major and minor scales by applying the tonal harmony
- Polyphonic techniques
  - horizontal: poly-melody: counter-point
  - vertical: chords, succession of harmony
- Major and minor scales
- Circle of fifth
  - 12 Major and minor keys



# What is pitch ?

# How pitches are organized ?

- **Chords:**

- set of notes considered as a whole in terms of harmony
  - most often, notes are played (almost) simultaneously
    - but also can be arpeggiated

- Chords on Major scale

I      II      III      IV      V      VI      VII      VIII

C MAJ<sup>7</sup>    D MIN<sup>7</sup>    E MIN<sup>7</sup>    F MAJ<sup>7</sup>    G<sup>7</sup>    A MIN<sup>7</sup>    B MIN<sup>7(5)</sup>    C MAJ<sup>7</sup>

## Chord dictionary

the chord

Root	Major	m	+	6	m6	7	m7	ma7	=	9
C	C	Cm	C+	C6	Cm6	C7	Cm7	Cma7	C°	C9
C# D	D	D#	C#m	D+6	D#m6	D#7	C#m7	D#ma7	C#°	D#9
D	D	Dm	D+	D6	Dm6	D7	Dm7	Dma7	D°	D9
D# E	E	E#	Em	E+6	E#m6	E#7	Em7	Em#7	E#°	E#9
E7	Em7	Em#7							E°	E9
F7	Fm7	Fma7							F°	F9
F#7	F#m7	F#ma7							F#°	F#9
G7	Gm7	Gma7							G°	G9
A#7	G#m7	A#ma7							A#°	A#9
A7	Am7	Am#7							A°	A9
B67	B#m7	B#ma7							B#°	B#9
B7	Bm7	Bma7							B°	B9

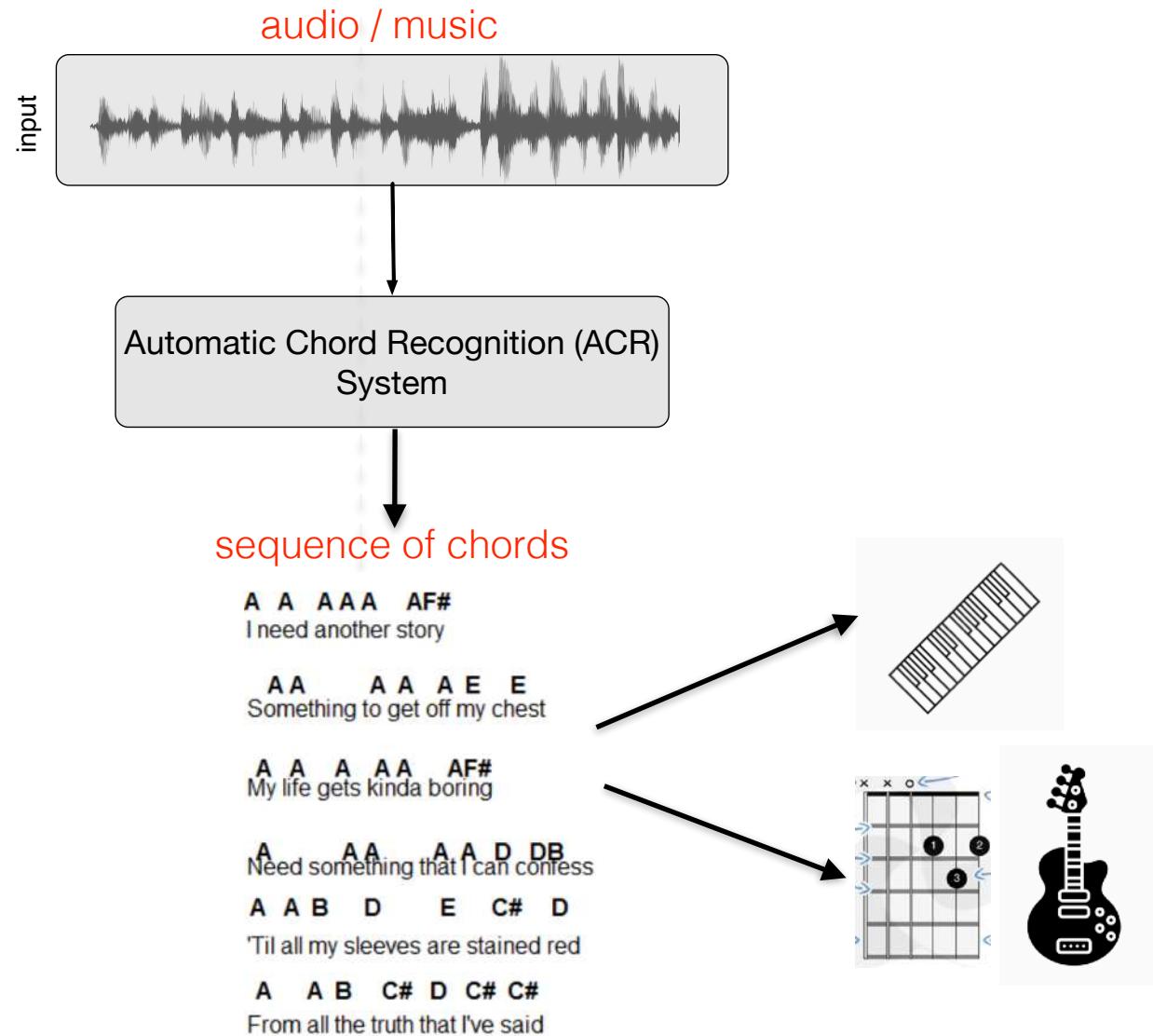
major 7 (C-M7: c, e, g, b), minor 7 (C-m7: c, e $\flat$ , g, b), sus-  
 (C-sus2: c, d, g / C-sus4: c, f, g), augmented (C-aug:  
 diminished (C-dim: c, e $\flat$ , g $\flat$ )  
 major 7 (C-M7: c, e, g, b), minor 7 (C-m7: c, e $\flat$ , g,  
 minor 7 (C-7: c, e, g, b $\flat$ ), major 6 (C-M6: c, e, g, a),  
 (C-m6: c, e $\flat$ , g, a) ...  
 major 9 (C-M9: c, e, g, b, d), dominant 9 (C-9: c,  
 d $\sharp$ )

- Most common chord sequences

- $V \rightarrow I$
  - $II \rightarrow V \rightarrow I$
  - Anatole
  - The magic 4-chords succession in pop-music
    - $I \rightarrow V \rightarrow VI \rightarrow IV$
    - $VI \rightarrow IV \rightarrow I \rightarrow V$

# Automatic Chord Recognition (ACR)

# Automatic Chord Recognition (ACR)



# Automatic Chord Recognition (ACR)

## Chord dictionary, chord vocabulary

### – En-harmonicity:

- we suppose that C# is the same as Db

### – Different chord **types**:

- triads (3 notes)
  - Major, minor, diminished, augmented
- tetrads (4 notes)
  - Seventh, Sixth
- extended chords
  - ...

### – Missing notes

### – Chord **inversions**

- which one is the bass note ?
  - $(1,3,5)/1 \rightarrow (1,3,5)/3$

Chord Type	Shorthand Notation	Components List
Triad Chords:		
Major	maj	(3,5)
Minor	min	(b3,5)
Diminished	dim	(b3,b5)
Augmented	aug	(3,#5)
Seventh Chords:		
Major Seventh	maj7	(3,5,7)
Minor Seventh	min7	(b3,5,b7)
Seventh	7	(3,5,b7)
Diminished Seventh	dim7	(b3,b5,bb7)
Half Diminished Seventh	hdim7	(b3,b5,b7)
Minor (Major Seventh)	minmaj7	(b3,5,7)
Sixth Chords:		
Major Sixth	maj6	(3,5,6)
Minor Sixth	min6	(b3,5,6)
Extended Chords:		
Ninth	9	(3,5,b7,9)
Major Ninth	maj9	(3,5,7,9)
Minor Ninth	min9	(b3,5,b7,9)
Suspended Chords:		
Suspended 4th	sus4	(4,5)

$$C:min7 \equiv C:(b3,5,b7)$$

$$C:min7(*5,11) \equiv C:(b3,b7,11)$$

$$C \equiv C:maj \equiv C:(3,5)$$

$$A/3 \equiv A:maj/3 \equiv A:(3,5)/3$$

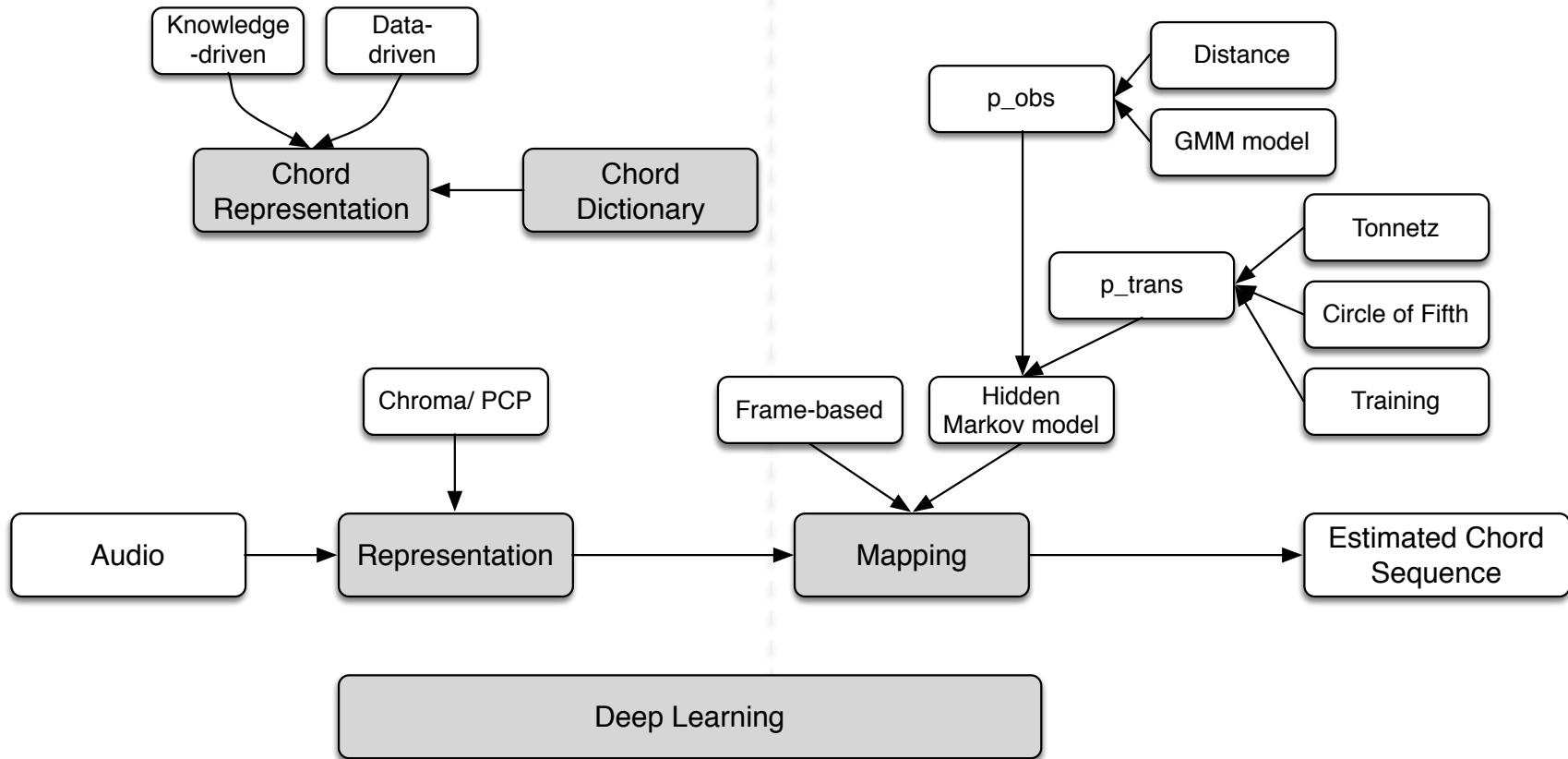
$$C:maj(4) \equiv C:(3,4,5)$$



# Automatic Chord Recognition (ACR) Systems

# Automatic Chord Recognition (ACR) Systems

## Automatic Chord Recognition System



# Automatic Chord Recognition (ACR) Systems

## Brief overview of system evolution

- as usual, the **first systems** define the task, the performance measures, and provide a first test-set; **later systems** deals with scalability issues and create large test-set; **current systems** use this large dataset to train systems using deep-learning
- Frame-based/ template-based approach
  - **1999** → T. Fujishima. "Realtime chord recognition of musical sound: a system using common lisp music". In Proc. of ICMC, 1999.
- Hidden-Markov-Model (HMM) based approaches
  - **2003** → A. Sheh and D. P. W. Ellis. "Chord segmentation and recognition using em-trained hidden Markov models". In Proc. of ISMIR, 2003
  - **2007** → H. Papadopoulos and G. Peeters. "Large-scale study of chord estimation algorithms based on chroma representation". In Proc. of IEEE CBMI, 2007
- Splitting into bass/middle/chroma
  - **2012** → Yizhao Ni, Matt McVicar, Raul Santos-Rodriguez, and Tijl De Bie. "An end-to-end machine learning system for harmonic analysis of music". IEEE TASLP, 2012.
- Deep learning approaches
  - **2013** → Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent. "Audio chord recognition with recurrent neural networks". In ISMIR, 2013
  - **2016** → Filip Korzeniowski and Gerhard Widmer. "Feature learning for chord recognition: the deep chroma extractor". In ISMIR, 2016.
  - **2017** → B. McFee and J. P. Bello. "Structured training for large-vocabulary chord recognition". In Proc. of ISMIR, 2017
  - **2021** → C. Weiß and G. Peeters. "Training deep pitch-class representations with a multi-label CTC loss". In Proc. of ISMIR, 2021



# Automatic Chord Recognition (ACR) Systems

## (1) Frame-based/ template-based approach

### – Frame-based:

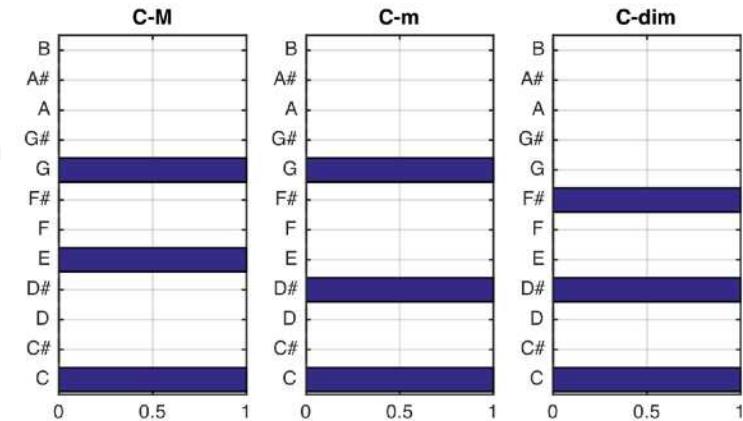
- chords are assigned to each temporal frame independently of each others

### – Template-based:

- chords are assigned by comparing chord-templates to chroma vectors

### – Chord templates

- 12 dimensional vector (if chroma has 12 dimensions) with value of
  - 1 if chroma exists in the chord
  - 0 if chroma does not exist in the chord
- Template  $G_a(c)$ 
  - chord-name,  $a \in \{C\text{-M}, C\text{-m}, C\#\text{-M}, C\#\text{-m}, \dots\}$
  - chroma-index,  $c \in [0, 12[$
- Variations [Gomez, 2006]
  - profiles can be extended to the audio case (harmonics of each pitch) by considering a contribution of the  $h$  harmonic of each pitch with an amplitude of  $0.6^{h-1}$ :  
use the first  $H=4$  harmonics



[E. Gomez. Tonal description of polyphonic audio for music content processing. *INFORMS Journal on Computing*, 2006.]

# Automatic Chord Recognition (ACR) Systems

## (1) Frame-based/ template-based approach

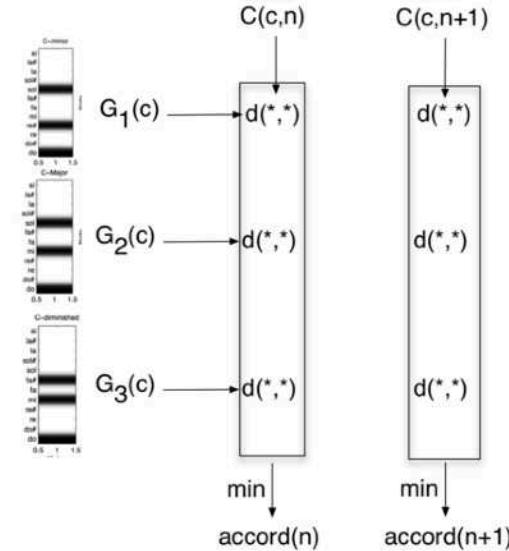
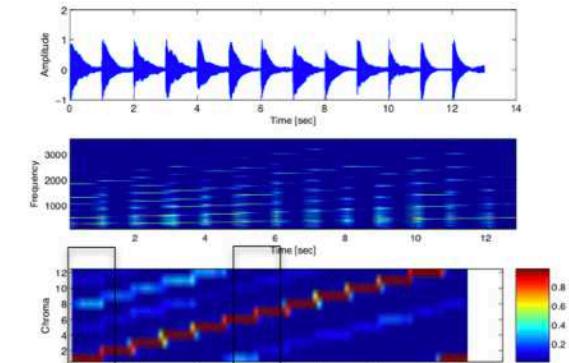
### – Frame-based estimation

### – Inputs:

- $C(c, n)$ : chroma vector at time  $n$
- $G_a(c) \quad a \in \{C\text{-M}, C\text{-m}, C\#\text{-M}, C\#\text{-m}, \dots\}$ : set of chord -templates

### – Compute

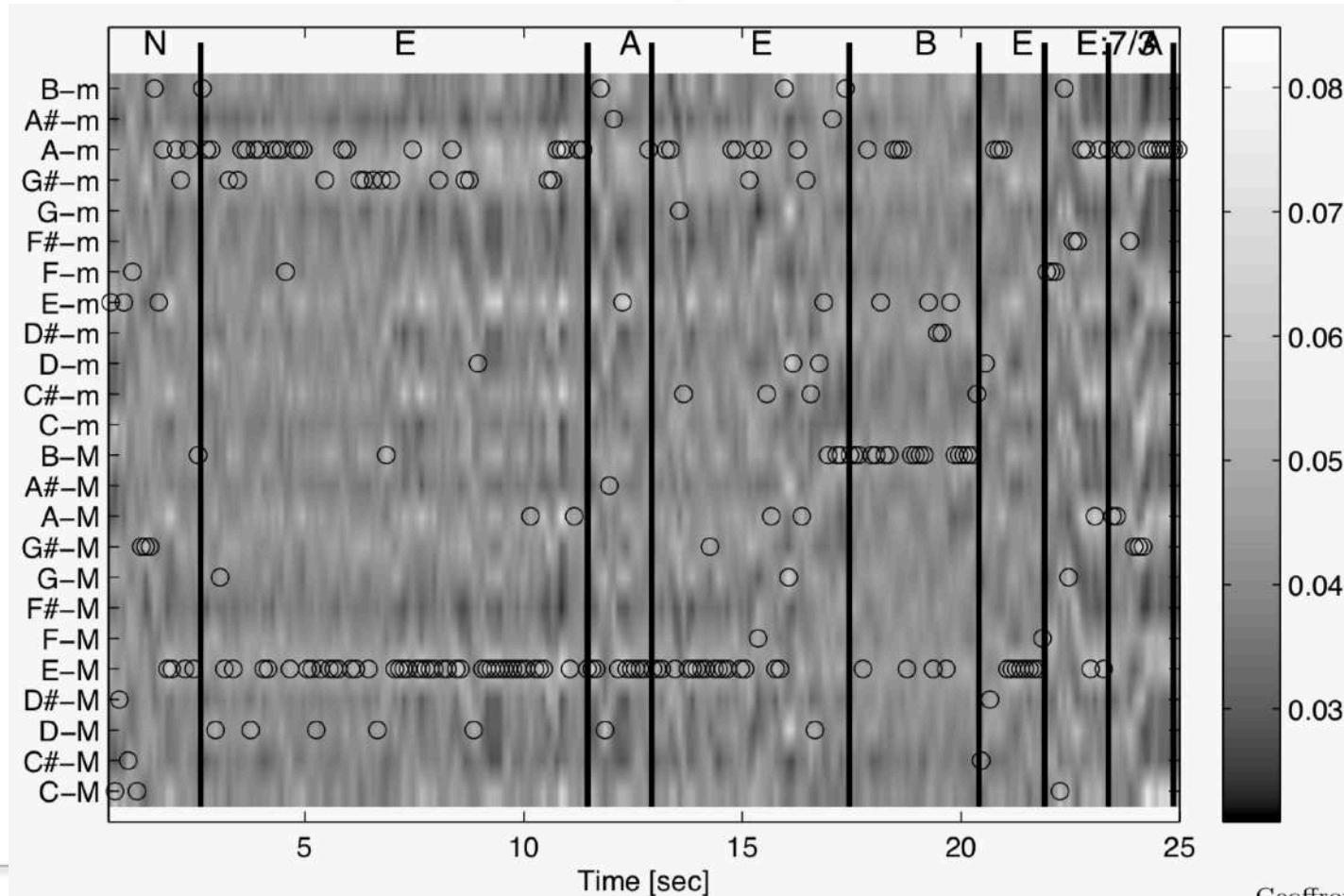
- $d(C(c, n), G_a(c))$  with  $d$ 
  - Euclidean distance
  - cosine distance
- Choose the chord with the smallest distance
  - $\arg \min_a d(C(c, n), G_a(c))$



# Automatic Chord Recognition (ACR) Systems

## (1) Frame-based/ template-based approach

– Result example



# Automatic Chord Recognition (ACR) Systems

## (2) HMM-based approach

### Chord Segmentation and Recognition using EM-Trained Hidden Markov Models

Alexander Sheh and Daniel P.W. Ellis  
LabROSA, Dept. of Electrical Engineering,  
Columbia University, New York, NY 10027 USA  
[asheh79, dpwe]@ee.columbia.edu

#### Abstract

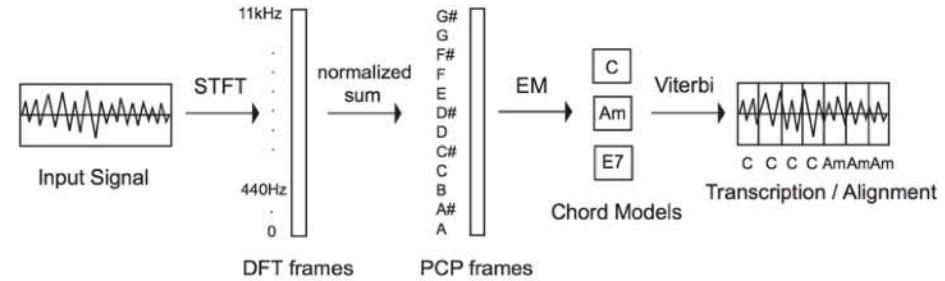
Automatic extraction of content description from commercial audio recordings has a number of important applications, from indexing and retrieval through to novel musicological analyses based on very large corpora of recorded performances. Chord sequences are a description that captures much of the character of a piece in a compact form and using a modest lexicon. Chords also have the attractive property that a piece of music can (mostly) be segmented into time intervals that consist of a single chord, much as recorded speech can (mostly) be segmented into time intervals that correspond to words. In this work, we build a system for automatic chord transcription using speech recognition tools. For features we use "pitch class profile" vectors to emphasize the tonal content of the signal, and we show that these features far outperform cepstral coefficients for our task. Sequence recognition is accomplished with hidden Markov models (HMMs), directly analogous to subword models in a speech recognizer, and trained by the same Expectation-Maximization (EM) algorithm. Crucially, this allows us to use as input only the chord sequences for our training example, without requiring the training set to contain the chords — which are determined automatically during training. Our results on a small set of 20 early Beatles songs show frame-level accuracy of around 75% on a forced-alignment task.

**Keywords:** audio, music, chords, HMM, EM.

#### 1 Introduction

The human auditory system is capable of extracting rich and meaningful data from complex audio signals. Machine listening research attempts to model this process using computers. In the music domain, there has been limited success when the input signal or analysis is relatively simple, i.e. single instrument,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. ©2003 Johns Hopkins University.



[A. Sheh and D. P. W. Ellis. "Chord segmentation and recognition using EM-trained hidden Markov models". In Proc. of ISMIR, 2003]

[H. Papadopoulos and G. Peeters. Large-scale study of chord estimation algorithms based on chroma representation. In Proc. of IEEE CBMI, 2007]

# Automatic Chord Recognition Systems

## (2) HMM-based approach

### – Defining HMM for chord recognition

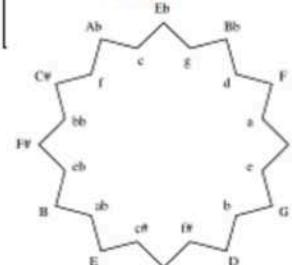
- **Observations  $O_t$** 
  - sequence of chroma/PCP vectors over time  $t$
- **States  $S_1, S_2, \dots, S_N$** 
  - the chord labels (C-M, C-m, C#-M, ...)
- **Initial state distribution  $\pi = \{\pi_i\}$** 
  - set to uniform distribution (if we don't have information)
- **Emission probabilities  $b_j(O_t) = p(O_t | q_t = S_j)$**
- **Transition probability  $a_{ij} = P(q_t = S_j | q_{t-1} = S_i)$**

### – Goal:

- Given a sequence of observations  $O = O_1, O_2, \dots, O_T$  and a model  $\lambda = \{A, B, \pi\}$ 
  - find the most likely  $Q = q_1, q_2, \dots, q_T \rightarrow$  **Viterbi decoding algorithm**

### Hidden Markov Model

#### Modeling transition between



FM

em

CM

am

GM

#### Modeling chords



# Automatic Chord Recognition Systems

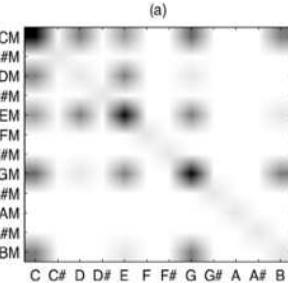
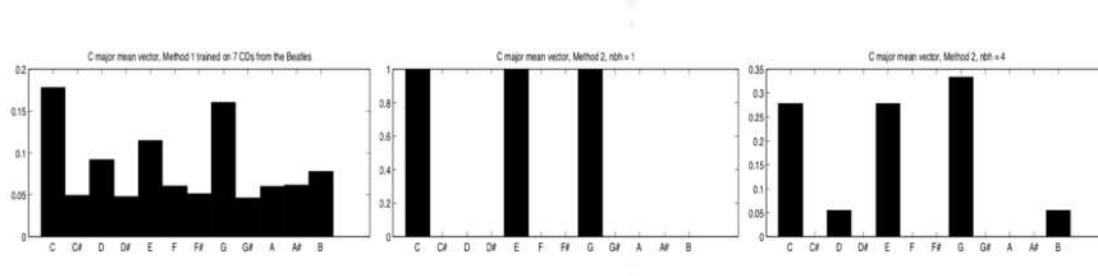
## (2) HMM-based approach

### – Emission probabilities $b_j(O_t) = p(O_t | q_t = S_i)$

- can be **trained** on a corpus using Baum-Welch algorithm
- can be **based** on the (normalised) distance to chord-templates (see before)
- can be based on manually-tuned statistical models
  - multivariate Gaussian models with parameters  $\mu$  and  $\Sigma$
  - mean vectors and covariance matrices reflect musical knowledge

mean vectors are 12-dim vectors with 1 if the note belongs to the chord

considers the correlation between the chroma vectors corresponding to the pitch of the notes belonging to a given chord. In our method, we also consider the correlation between the harmonics of each note



# Automatic Chord Recognition Systems

## (2) HMM-based approach

- Transition probability between chords:
  - can be **trained** using the Baum-Welch algorithm
  - can be **trained** on a symbolic-music corpus
    - Guitar Tab, Real Book

[Intro]  
C Am C Am

[Verse 1]  
C Am  
I heard there was a secret chord  
C Am  
That David played and it pleased the lord  
F G C G  
But you don't really care for music, do you?  
C F G  
Well it goes like this the fourth, the fifth  
Am F  
The minor fall and the major lift  
G E7 Am  
The baffled king composing hallelujah

[Chorus]  
F Am F C G C Am C Am  
Hallelujah, hallelujah, hallelujah, hallelu-u-u-jah . . .

Chord Progression:

B<sup>b</sup> | F<sub>7</sub> | C<sub>-7</sub> F<sub>7</sub> |  
B<sup>b</sup> | B<sub>07</sub> | F<sub>7</sub> | A<sub>-7</sub> D<sub>7</sub> |  
G<sub>-7</sub> | C<sub>7</sub> | F<sub>7</sub> D<sub>7</sub> | G<sub>-7</sub> C<sub>7</sub> |  
B<sup>b</sup> | F<sub>7</sub> | B<sup>b</sup> | C<sub>-7</sub> F<sub>7</sub> |  
B<sup>b</sup> | B<sub>07</sub> | F<sub>7</sub> | A<sub>-7</sub> D<sub>7</sub> |  
G<sub>-7</sub> | C<sub>7</sub> | F<sub>7</sub> D<sub>7</sub> | G<sub>-7</sub> C<sub>7</sub> |

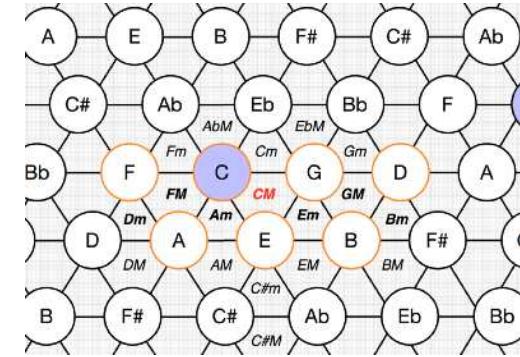
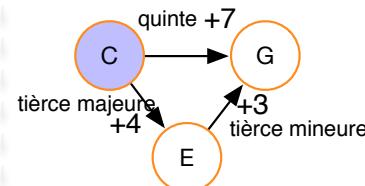
# Automatic Chord Recognition Systems

## (2) HMM-based approach

### – Transition probability between chords:

- can be **based on musical rules**

- distance between chords in Tonnetz space [Euler, 1739]



- distance between chords in the Circle-of-Fifth

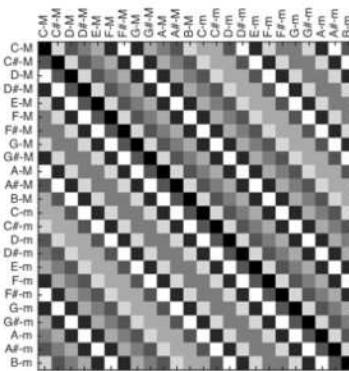
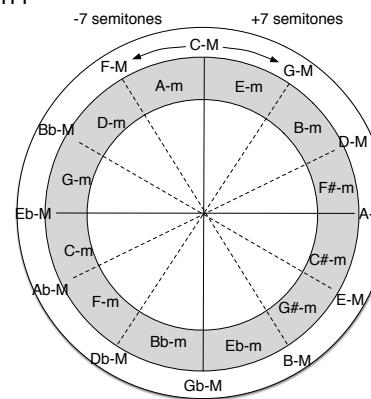
sequence of fifths

relative Major-minor

Example:

G-M to C-M (consonance),

G-M to Db-M (dissonance)



# Automatic Chord Recognition Systems

## (2) HMM-based approach

- **Transition probability between chords:**

- can be **based on musical rules**

- represent prototypical chord progressions

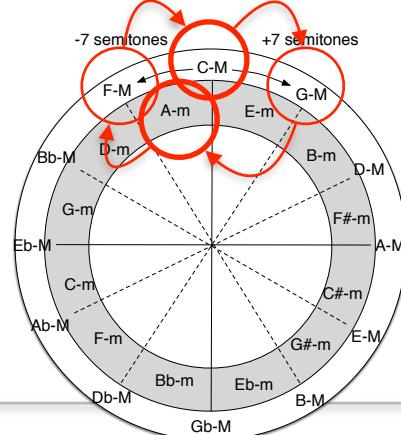
$V \rightarrow I:$

$II \rightarrow V \rightarrow I:$

Anatole ( $VI \rightarrow II \rightarrow V \rightarrow I$ ):

12-bar Blues:

Magic 4-chords sequence in pop-music:



$G7 \rightarrow Cmaj7$

$D-7 \rightarrow G7 \rightarrow Cmaj7$

$A-7 \rightarrow D-7 \rightarrow G7 \rightarrow Cmaj7$

$C7 \rightarrow F7 \rightarrow C7 \rightarrow C7 ||$

$F7 \rightarrow F7 \rightarrow C7 \rightarrow C7 ||$

$G7 \rightarrow F7 \rightarrow C7 \rightarrow G7 ||$

$C \rightarrow G \rightarrow A- \rightarrow F$

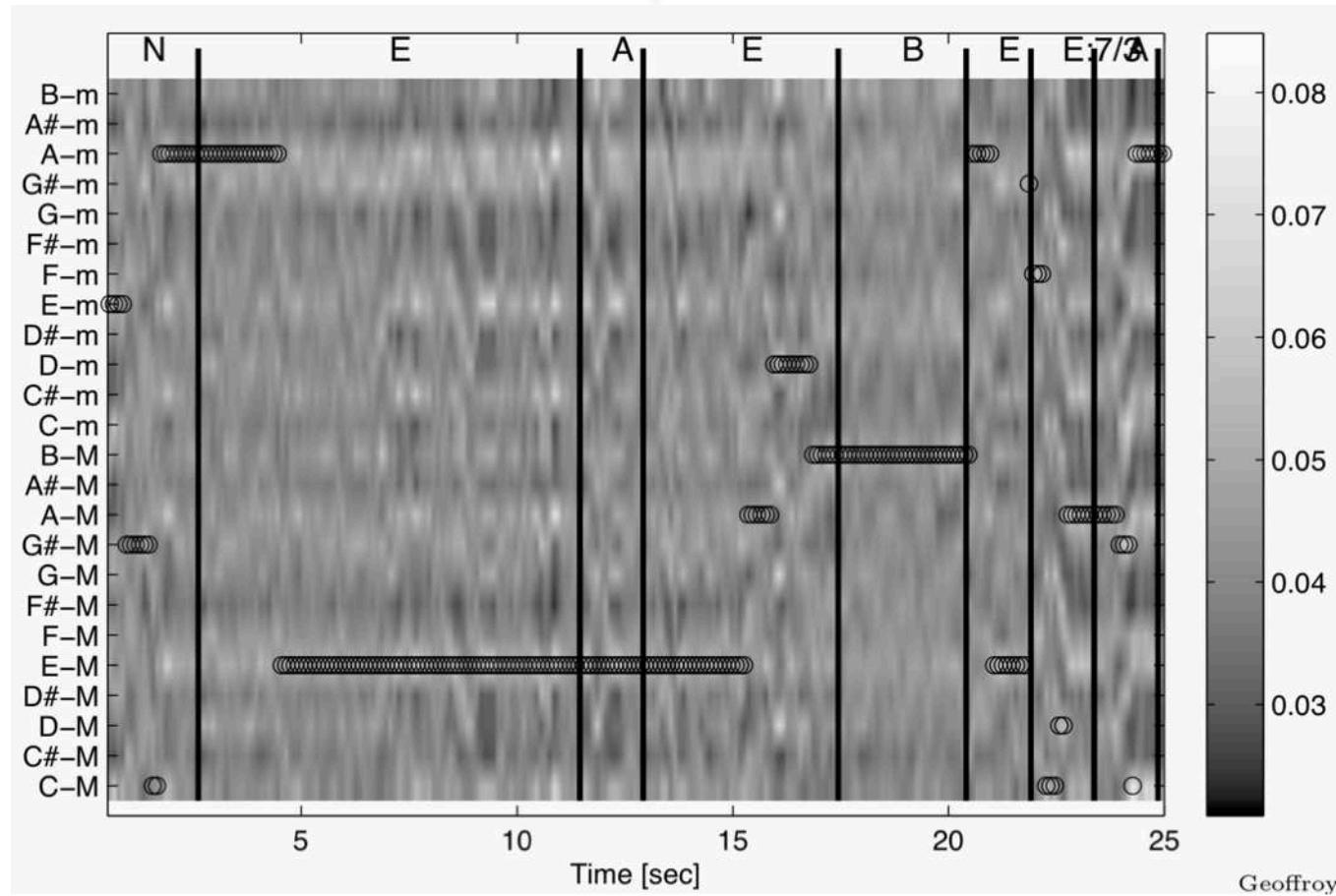
$A- \rightarrow F \rightarrow C \rightarrow G$



# Automatic Chord Recognition Systems

## (2) HMM-based approach

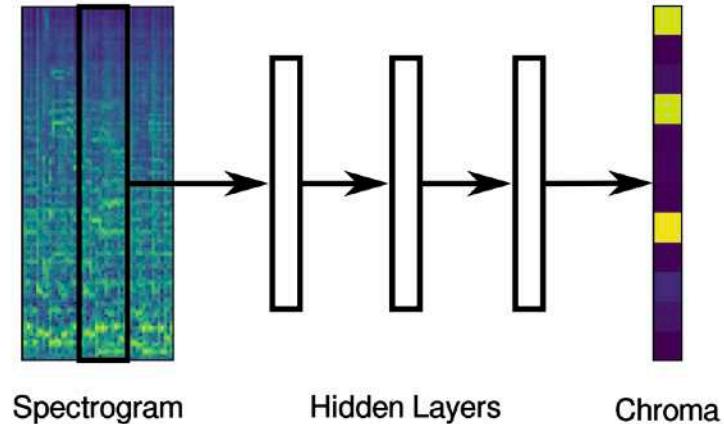
- Result example



# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach A

- **Goal:**
  - standard Chroma extractors = too noisy features
  - replace the Chroma front-end by learned features
    - encode harmonic information important for chord recognition, while being robust to irrelevant interferences
    - train a 3-layers MLP to output a ground-truth chroma representation
    - ground-truth ? Chroma corresponding to the notes of the chord)
    - feeding the network with an audio spectrum with context instead of a single frame as input
  - **Deep Chroma**



	Btls	Iso	RWC	RW	Total
$C$	$71.0 \pm 0.1$	$69.5 \pm 0.1$	$67.4 \pm 0.2$	$71.1 \pm 0.1$	$69.2 \pm 0.1$
$C_{Log}^W$	$76.0 \pm 0.1$	$74.2 \pm 0.1$	$70.3 \pm 0.3$	$74.4 \pm 0.2$	$73.0 \pm 0.1$
$S_{Log}$	$78.0 \pm 0.2$	$76.5 \pm 0.2$	$74.4 \pm 0.4$	$77.8 \pm 0.4$	$76.1 \pm 0.2$
$C_D$	$80.2 \pm 0.1$	$79.3 \pm 0.1$	$77.3 \pm 0.1$	$80.1 \pm 0.1$	$78.8 \pm 0.1$

$C$ : standard chroma from CQT

$C_{Log}^W$ : chromagram with frequency weighting and logarithmic compression

$S_{Log}$ : quarter-tone spectrogram

$C_D$ : deep-chroma

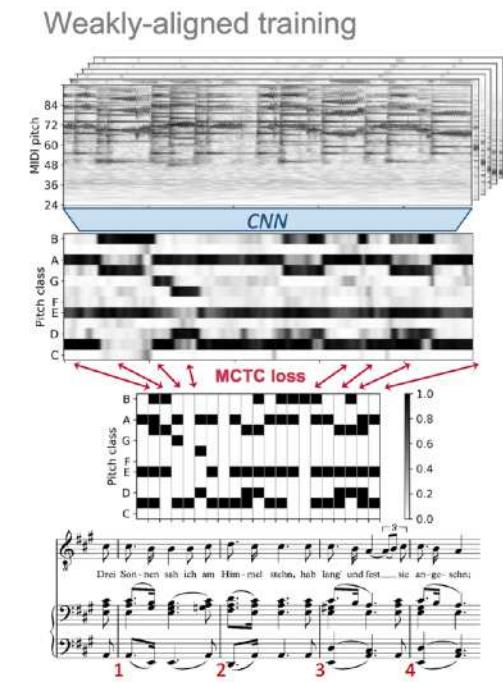
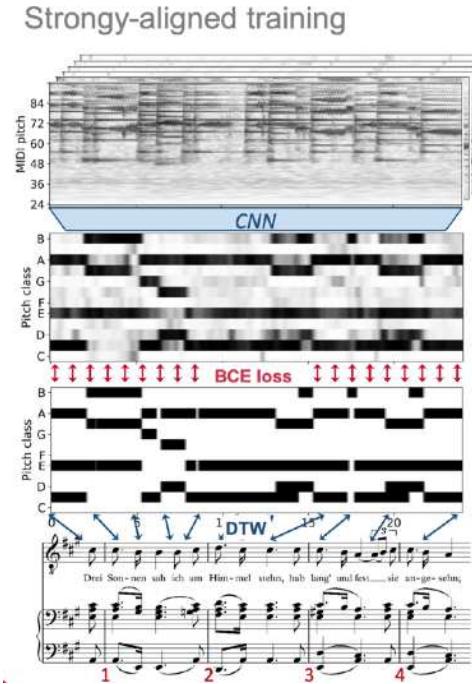


# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach A

### – Goal:

- replace the Chroma/PCP front-end by learned features
- Ground-truth ?
  - Aligned pitches (costly)
  - Non-aligned pitches (CTC)



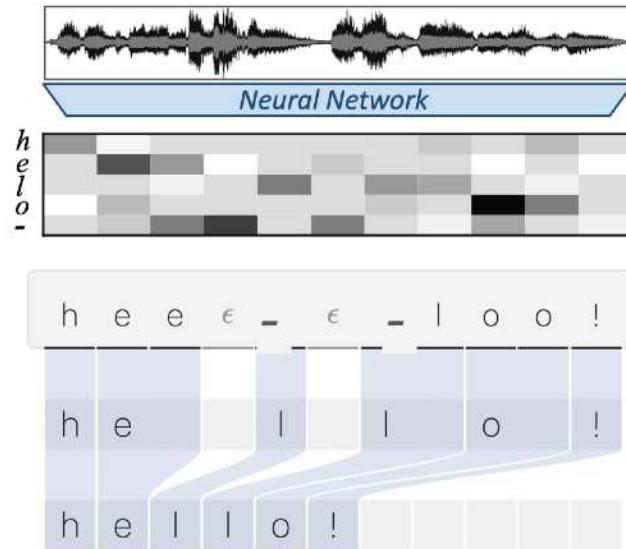
[C. Weiß and G. Peeters. "Training deep pitch-class representations with a multi-label CTC loss". In Proc. of ISMIR, 2021] [LINK](#)

# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach A

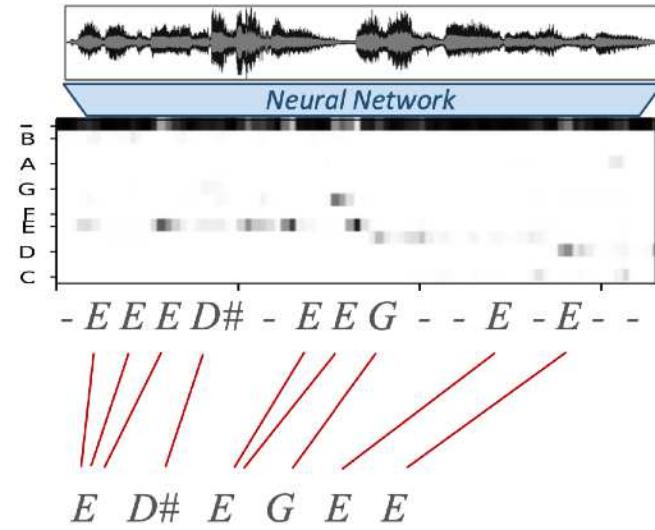
### – Connectionist Temporal Classification (CTC) Loss

Automatic Speech Recognition



Graves, Fernández, Gomez, Schmidhuber:  
Connectionist temporal classification:  
labelling unsegmented sequence data with  
recurrent neural networks. Proc. ICML 2006

Monophonic pitch-class estimation



Zalkow, Müller: Using Weakly Aligned  
Score-Audio Pairs to Train Deep Chroma  
Models for Cross-Modal Music Retrieval.  
Proc. ISMIR 2020

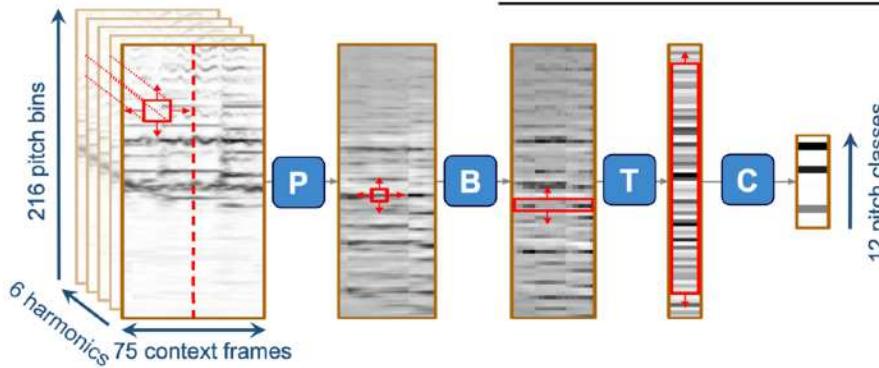


# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach A

### – CNN Architecture

- Input: **Harmonic-CQT**
- Simple **5-layerCNN**
- Roughly **48k parameters**
- Pre-filtering, **B**inning to midi-pitches ( $216 \rightarrow 72$ ), **T**emporal reduction ( $75 \rightarrow 1$ ), **C**hroma reduction ( $72 \rightarrow 12$ )
- **Input:** Harmonic CQT



Layer	Kernel size	Output shape	# Parameters
Layer norm.		$(T+74, 216, 6)$	2592
P Conv2D, MaxPool	$15 \times 15$	$(T+74, 216, 20)$	27020
B Conv2D, MaxPool	$3 \times 3$	$(T+74, 72, 20)$	3620
T Conv2D	$75 \times 1$	$(T, 72, 10)$	15010
Conv2D	$1 \times 1$	$(T, 72, 1)$	11
C Conv2D	$1 \times 61$	$(T, 12+P, Q)$	$Q(62+73 \cdot P)$
<b>Total</b>			48253 $+Q(62+73 \cdot P)$

[C. Weiß and G. Peeters. "Training deep pitch-class representations with a multi-label CTC loss". In Proc. of ISMIR, 2021] [LINK](#)

# Automatic Chord Recognition (ACR) Systems

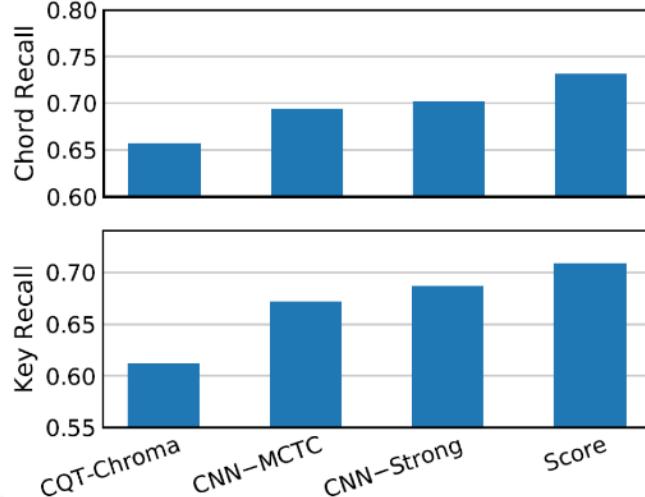
## (2) Deep-learning-based approach A

### – Evaluation

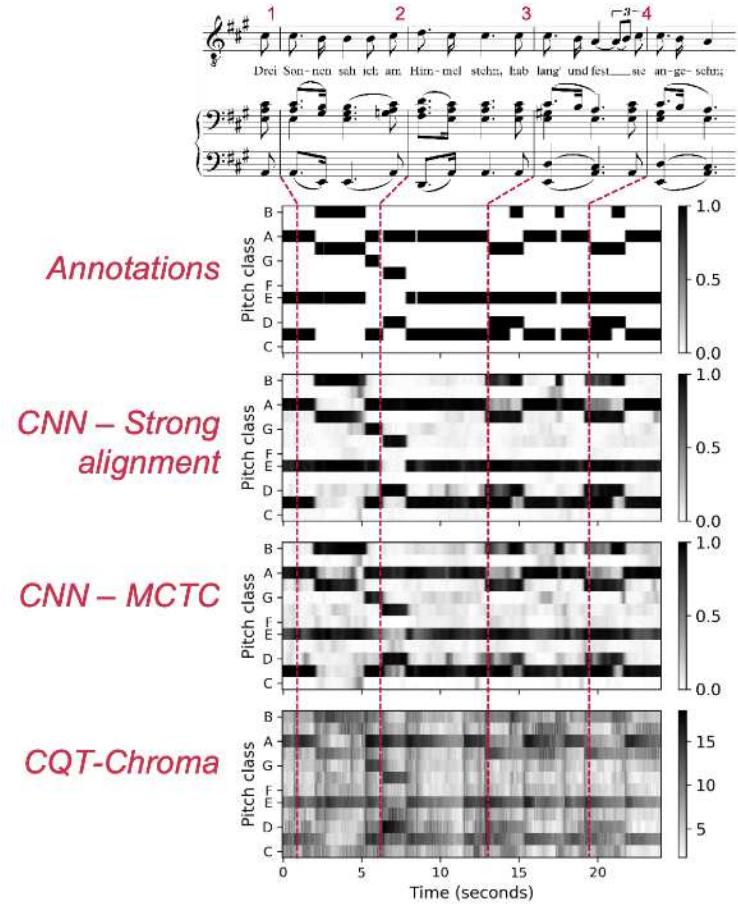
- Cosine similarity (CS), Average precision (AP)

Model/Loss	P	R	F	CS	AP
All-Zero	0	0	0	0.486	0.211
CQT-Chroma	0.512	0.681	0.579	0.701	0.594
CNN – SCTC	<b>0.850</b>	0.048	0.090	0.520	0.416
CNN – MCTC:NE	0.747	0.775	0.758	0.802	0.798
CNN – MCTC:WE	0.762	<b>0.853</b>	0.802	0.830	0.851
CNN – Strong alignment	<b>0.850</b>	0.790	<b>0.818</b>	<b>0.860</b>	<b>0.886</b>

### – Application: Chord and Key estimation



### Application: Visualization



[C. Weiß and G. Peeters. "Training deep pitch-class representations with a multi-label CTC loss". In Proc. of ISMIR, 2021] [LINK](#)

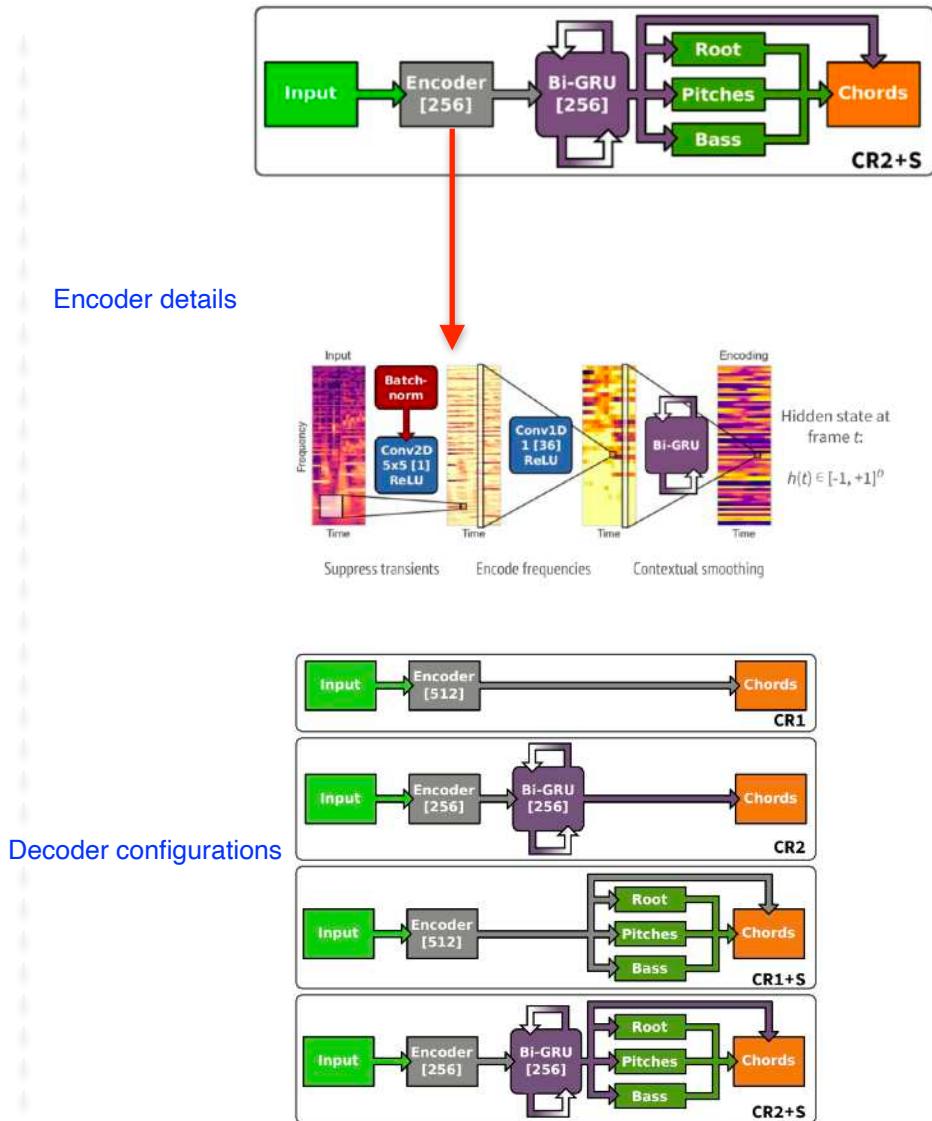
# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach B

- Goal 1:
  - End-to-end system

- Encoder:
  - Input:  $T \times F$  time-series of log-power constant-Q transform (CQT) spectra
  - First layer : can be interpreted as a harmonic saliency enhancer, as it tends to learn to suppress transients and vibrato while emphasizing sustained tones.
  - Second layer summarizes the pitch content of each frame, and can be interpreted as a local feature extractor

- Decoder:
  - 4 architectures



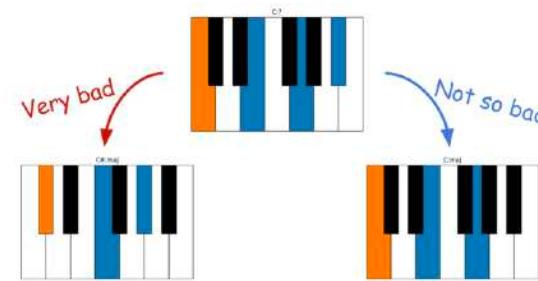
[B. McFee and J. P. Bello. "Structured training for large-vocabulary chord recognition". In Proc. of ISMIR, 2017] [LINK](#)

# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach B

- Goal 1:
  - End-to-end system
- Goal 2:
  - Large chord vocabularies
    - Classes are not well-separated  
 $C:7 = C:\text{maj} + m7$   
C:sus4 vs. F:sus2
    - Class distribution is non-uniform
    - Rare classes are hard to model
  - Take into account the fact that **some mistakes are better than others**

$$14 \times 12 + 2 = 170 \text{ classes}$$



[B. McFee and J. P. Bello. "Structured training for large-vocabulary chord recognition". In Proc. of ISMIR, 2017] [LINK](#)

# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach B

### Idea:

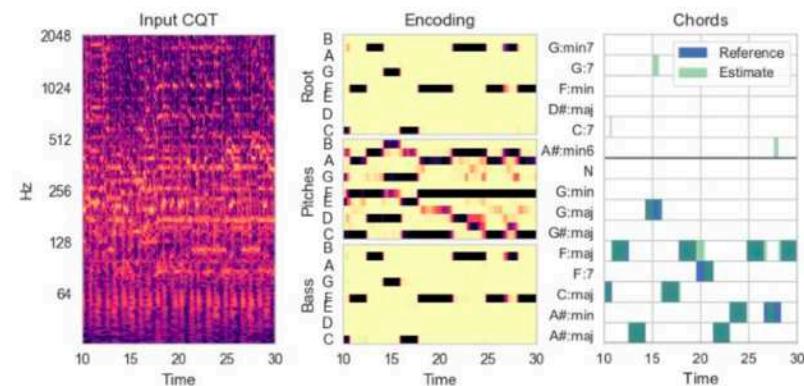
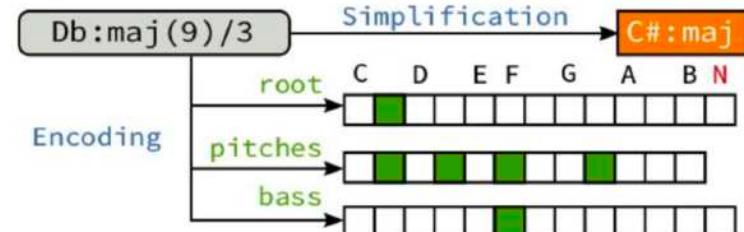
- Exploit the implicit chord space structure
- Represent chord labels as **binary encodings**
  - Similar chords with different labels will have similar encodings
  - Dissimilar chords will have dissimilar encodings

### Learning problem:

- Predict the **encoding** from audio
- Learn to decode into chord **labels**

### Model architecture

- Input: constant-Q spectral patches
- Per-frame outputs:
  - Root [multi-class, C=13]
  - Pitches [multi-label, C+12]
  - Bass [multi-class, C=13]
  - Chords [multi-class, C=170]
- Convolutional-recurrent architecture (encoder-decoder)
- End-to-end training



[B. McFee and J. P. Bello. "Structured training for large-vocabulary chord recognition". In Proc. of ISMIR, 2017] [LINK](#)

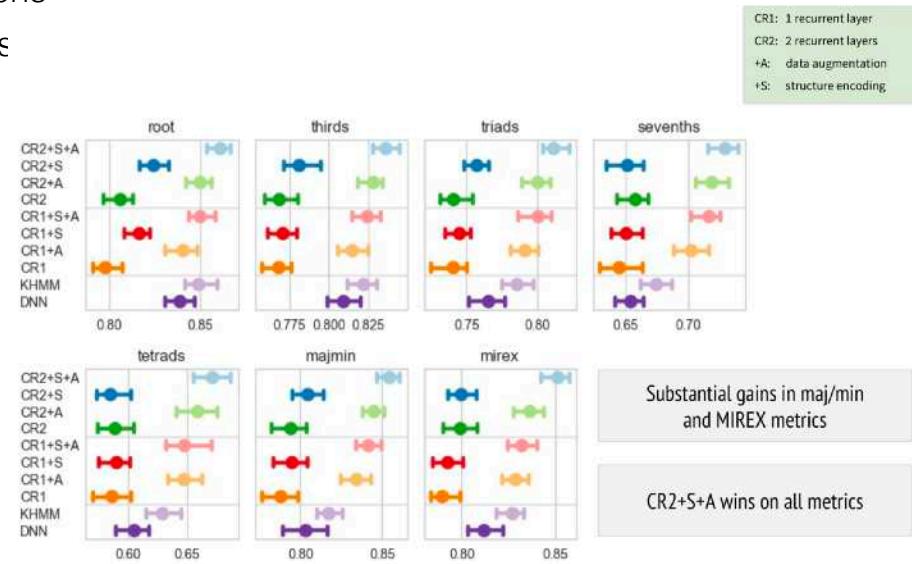
# Automatic Chord Recognition (ACR) Systems

## (2) Deep-learning-based approach B

- **What about root bias?**
  - Quality and root should be independent
  - But the data is inherently biased
- **Solution: data augmentation!**
  - Pitch-shift the audio and annotations simultaneously
  - Each training track → ± 6 semitone shifts
    - All qualities are observed in all root positions
    - All roots, pitches, and bass values are obs
- **8 configurations**
  - ± data augmentation
  - ± structured training
  - vs. 2 recurrent layers
- **1217 recordings**
  - (Billboard + Isophonics + MARL corpus)
  - 5-fold cross-validation
- **Baseline models:**
  - DNN [Humphrey & Bello, 2015]
  - KHMM [Cho, 2014]

### • Training details

- Keras / TensorFlow + pescador
- ADAM optimizer
- Early stopping @20, learning rate reduction @10
  - Determined by decoder loss
- 8 seconds per patch
- 32 patches per batch
- 1024 batches per epoch



[B. McFee and J. P. Bello. "Structured training for large-vocabulary chord recognition". In Proc. of ISMIR, 2017] [LINK](#)



# Automatic Chord Recognition (ACR)

## Evaluation

# Automatic Chord Recognition (ACR) Evaluation

## Task definition:

- [https://www.music-ir.org/mirex/wiki/2019:Audio\\_Chord\\_Estimation](https://www.music-ir.org/mirex/wiki/2019:Audio_Chord_Estimation)
- Given an audio track
  - estimate the set of temporal (start:end) segments and associated chord label

```
0.000000 2.612267 N
2.612267 11.459070 E
11.459070 12.921927 A
12.921927 17.443474 E
17.443474 20.410362 B
20.410362 21.908049 E
21.908049 23.370907 E:7/3
23.370907 24.856984 A
24.856984 26.343061 A:min/b3
26.343061 27.840748 E
27.840748 29.350045 B
29.350045 35.305963 E
35.305963 36.803650 A
36.803650 41.263102 E
41.263102 44.245646 B
44.245646 45.720113 E
45.720113 47.206190 E:7/3
47.206190 48.692267 A
48.692267 50.155124 A:min/b3
50.155124 51.652811 E
51.652811 53.138888 B
53.138888 56.111043 E
56.111043 65.131995 A
65.131995 68.150589 B
68.150589 71.192403 A
71.192403 74.199387 E
74.199387 75.697074 A
75.697074 80.236575 E
80.236575 83.208730 B
83.208730 86.221693 E
86.221693 87.736621 A
87.736621 89.257528 A:min/b3
89.257528 90.720385 E
```



# Automatic Chord Recognition (ACR) Evaluation

## Datasets:

- <https://www.audiocontentanalysis.org/data-sets/>
- QMUL Isophonics (Beatles, Carole King, Queen, Michael Jackson)
  - <http://isophonics.net/datasets>
- AIST RWC (Real World Computing)
  - <https://staff.aist.go.jp/m.goto/RWC-MDB/>
- McGill Billboard
  - [https://ddmal.music.mcgill.ca/research/The\\_McGill\\_Billboard\\_Project\\_\(Chord\\_Analysis\\_Dataset\)/](https://ddmal.music.mcgill.ca/research/The_McGill_Billboard_Project_(Chord_Analysis_Dataset)/)
- GuitarSet
  - <https://guitarset.weebly.com/>
- ...

# Automatic Chord Recognition (ACR) Evaluation

## Performance measures

### – Two criteria to evaluate

- **(1)** get the correct segment boundaries (independently of the labels)
- **(2)** get the correct label at each time

# Automatic Chord Recognition (ACR) Evaluation

## Example of results:

### Submissions

	Abstract	Contributors
CM1	<a href="#">PDF</a>	Chris Cannam  , Matthias Mauch 
JLCX1, JLCX2	<a href="#">PDF</a>	Junyan Jiang  , Ke Chen  , Wei Li  , Guangyu Xia 
SG1	<a href="#">PDF</a>	Franz Strasser  , Stefan Gaser 
FK2	<a href="#">[PDF]</a>	Florian Krebs  , Filip Korzeniowski  , Sebastian Böck 

### Results

#### Summary

All figures can be interpreted as percentages and range from 0 (worst) to 100 (best).

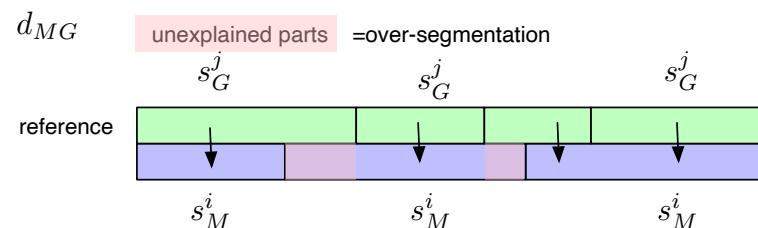
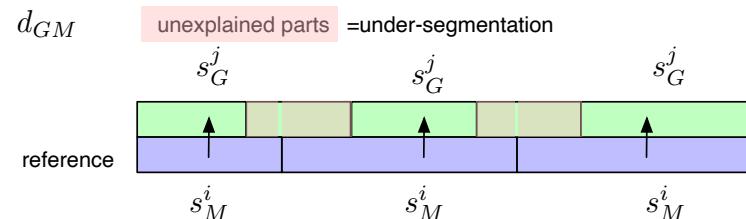
#### Isophonics2009

Algorithm	MirexRoot	MirexMajMin	MirexMajMinBass	MirexSevenths	MirexSeventhsBass	MeanSeg	UnderSeg	OverSeg
CM1	78.66	75.51	72.58	54.78	52.36	85.87	87.22	85.98
FK2	87.38	86.80	83.43	75.55	72.60	89.29	87.24	92.35
JLCX1	86.75	86.25	84.44	75.87	74.39	90.33	88.36	93.38
JLCX2	86.51	86.05	84.23	75.64	74.17	90.12	87.94	93.48
SG1	82.03	78.67	76.84	69.20	67.55	82.34	89.03	77.87

# Automatic Chord Recognition (ACR) Evaluation

## (1) get the correct segment boundaries (independently of the labels)

- Notation
  - ground-truth segments:  $S_G^j$
  - estimated (measured) segments:  $S_M^i$
- directional Hamming distance  $d_{GM}$ :
  - For each  $S_M^i$  find the segment  $S_G^j$  with the maximum overlap → then summing the differences
    - $d_{GM} = \sum_{S_M^i} \sum_{S_G^j \neq S_G^k} |S_M^i \cap S_G^j|$
    - where  $| . |$  denotes the length of a segment
    - $d_{GM}$  indicates missed boundaries (**under-segmentation**)
- inverse-directional Hamming distance  $d_{MG}$ :
  - For each  $S_G^j$  find the segment  $S_M^i$  with ...
    - $d_{MG}$  indicates segment fragmentation (**over-segmentation**)
- **Chord segmentation:**  $Q = 1 - \frac{\text{max. of } d_{GM} d_{MG}}{\text{total duration of song}}$



# Automatic Chord Recognition (ACR) Evaluation

## (2) get the correct label at each time

– **Chord symbol recall (CSR)**:  $CSR = \frac{\text{total duration of segments where annotation equal estimation}}{\text{total duration of annotated segments}}$

### – Need to choose a chord vocabularies

- Chord root-note only
- Major and minor: {N, maj, min}
- Seventh chords: {N, maj, min, maj7, min7, 7}
- Major and minor with inversions: {N, maj, min, maj/3, min/b3, maj/5, min/5}
- Seventh chords with inversions: {N, maj, min, maj7, min7, 7, maj/3, min/b3, maj7/3, min7/b3, 7/3, maj/5, min/5, maj7/5, min7/5, 7/5, maj7/7, min7/b7, 7/b7}

### – Vocabulary mapping

- G:7(#9) is mapped to G:maj because the interval set of G:maj, {1,3,5}, is a subset of the interval set of the G:7(#9), {1,3,5,b7,#9}
- G:7(#9) is mapped to G:7 instead because the interval set of G:7 {1, 3, 5, b7}
- non-possible mapping (excluded from evaluation)
  - D:aug or F:sus4(9) to {maj, min}

# Hidden Markov Model (HMM)

# Hidden Markov Model (HMM)

- Markov model
  - Andreï A. Markov (1856-1922): Russian mathematician
- Markov chain:
  - a stochastic process with discrete time  $t$ , each time can be in one of the possible discrete states  $S_{i \in \{1, \dots, N\}}$
  - $q_t$  is the value of the state at time  $t$
- First order Markov chain:
  - the prediction of the current state  $q_t$  depends only on the previous time  $q_{t-1}$

$$p(q_t | q_{t-1}, q_{t-2} \dots q_0) = p(q_t | q_{t-1})$$



source: [https://fr.wikipedia.org/wiki/Andre%C3%A9\\_Markov\\_\(math%C3%A9maticien\)](https://fr.wikipedia.org/wiki/Andre%C3%A9_Markov_(math%C3%A9maticien))

# Hidden Markov Model (HMM)

## Observed Markov model

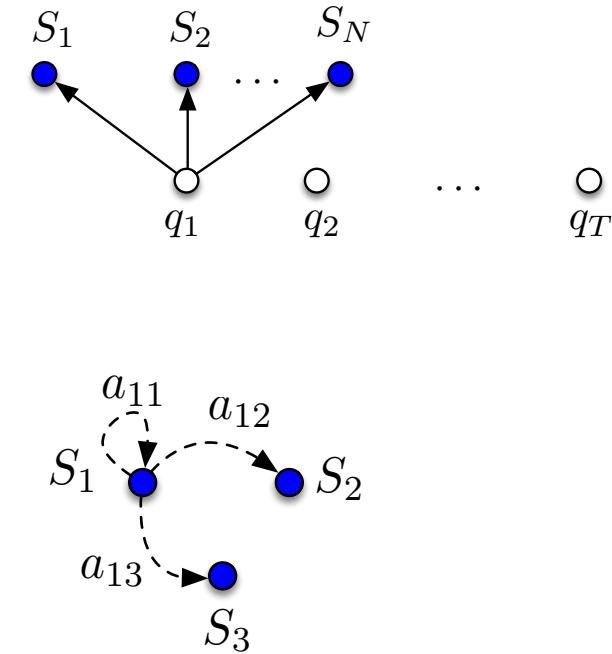
- The system is in one of a set of  $N$  distinct **states**  $S_1, S_2, \dots, S_N$
- We denote the **time** instants as  $t = 1, 2, \dots, T$
- We denote the **actual state** at time  $t$  as  $q_t$
- We denote the sequence of actual state  $Q = q_1, q_2, \dots, q_T$
- **Transition probability**

- Discrete first order Markov chain
  - $P(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots) = P(q_t = S_j | q_{t-1} = S_i)$
- is independent of time
  - $a_{ij} = P(q_t = S_j | q_{t-1} = S_i) \quad 1 \leq i, j \leq N$

with  $a_{ij} \geq 0$

$$\sum_j a_{ij} = 1$$

- **Initial state distribution**  $\pi = \{\pi_i\}$  with
  - $\pi_i = P(q_1 = S_i)$

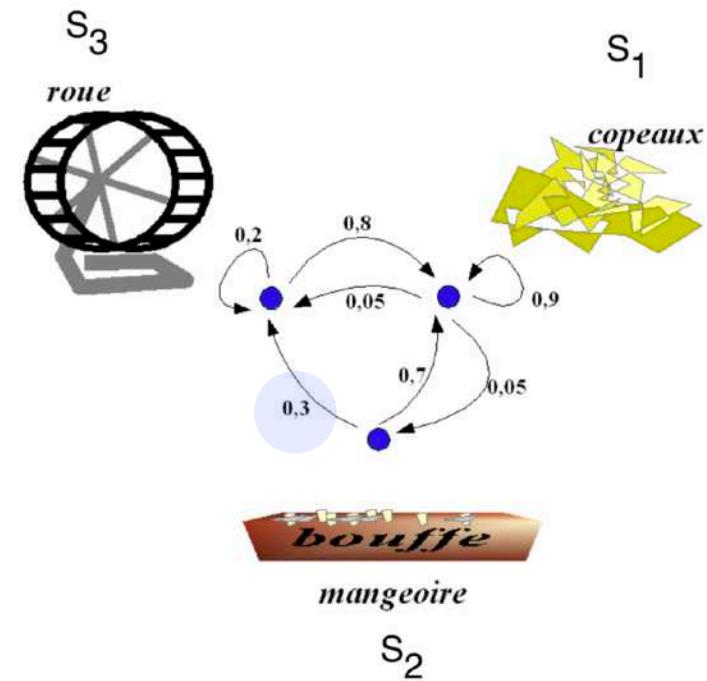


# Hidden Markov Model (HMM)

## Observed Markov model example

- Doudou the hamster has 3 states during a day
  - he sleeps:
    - he is in state  $S_1$  (sawdust shavings)
  - he eats:
    - he is in state  $S_2$  (feeder)
  - he does exercice:
    - he is in state  $S_3$  (hamster wheel)
- We can represent the succession of states a transition matrix  $A = (a_{ij})$  between states

$$A = \begin{pmatrix} 0.9 & 0.05 & 0.05 \\ 0.7^{(i=2,j=1)} & 0 & 0.3 \\ 0.8 & 0 & 0.2 \end{pmatrix} \quad \sum_j a_{ij} = 1$$



Modèle de Markov d'une journée de Doudou le hamster

source: [https://fr.wikipedia.org/wiki/Cha%C3%ABne\\_de\\_Markov](https://fr.wikipedia.org/wiki/Cha%C3%ABne_de_Markov)

# Hidden Markov Model (HMM)

## Hidden Markov model

- We do not observe directly the actual states  $Q = q_1, q_2, \dots, q_T$  but an emission of those
  - a sequence of **observations**  $O = O_1, O_2, \dots, O_T$

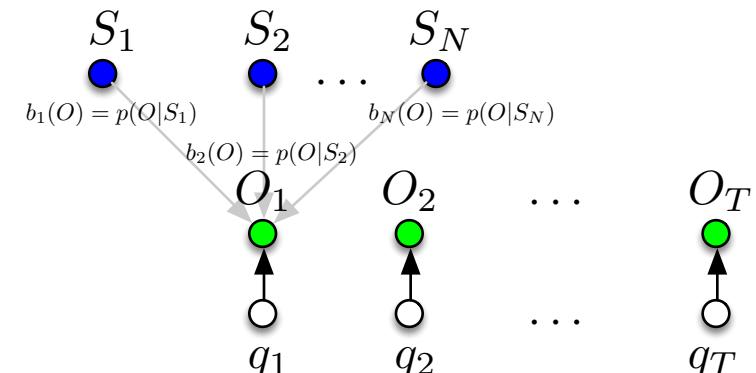
### Emission probability

- The observation is a probabilistic function of the state  $S_j$
- $b_j(O) = P(O | S_j)$

- **Discrete Observations:** set of symbols  $v_k$ 
  - $B = \{ b_j(k) \}$  with  $b_j(k) = P(O_t = v_k | q_t = S_j)$   
the probability of observing symbol  $k$  in state  $j$

- **Continuous Observation Densities:** Gaussian Mixture Model

$$b_j(O) = \sum_m c_{j(m)} \mathcal{N}(O | \mu_{j(m)}, \Sigma_{j(m)})$$



- We denote by  $\lambda = \{A, B, \pi\}$  the set of elements defining an HMM

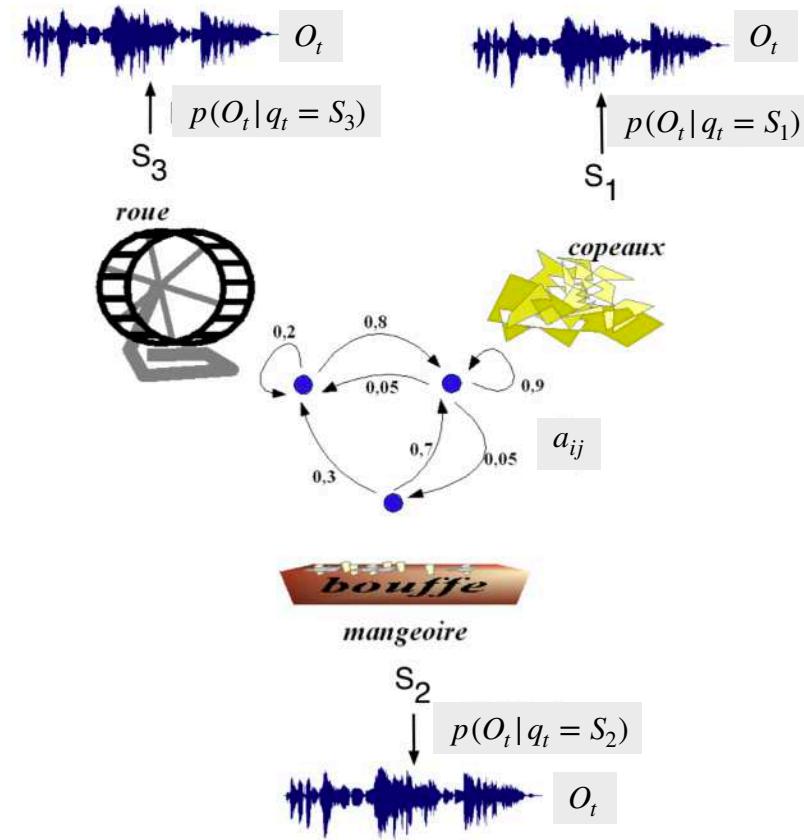
# Hidden Markov Model (HMM)

## Hidden Markov model example

- In a hidden Markov model, we do not observe directly the states  $Q = q_1, q_2, \dots, Q_T$ ,
  - they are "hidden"
- we observe an emission  $O_t$  of the states  $q_t$ 
  - example: we observe the sound  $O = O_1, O_2, \dots, O_T$  made by Doudou the hamster
- For each state, we can define
  - an emission probability given state  $S_j$ :

$$b_j(O_t) = p(O_t | q_t = S_j)$$

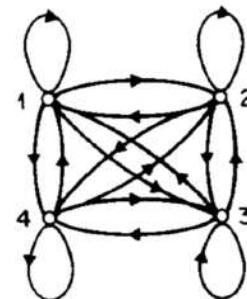
state: chord  
observation: chroma



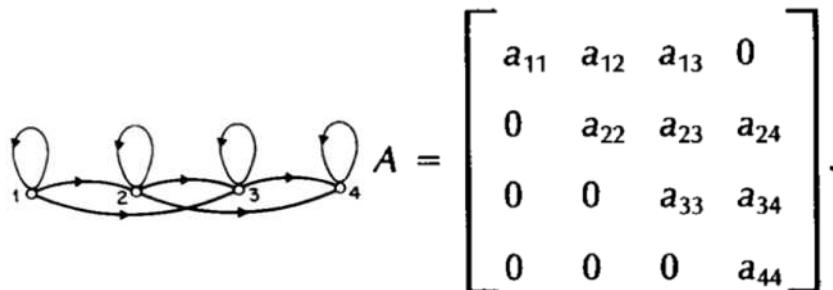
# Hidden Markov Model (HMM)

## Various topologies of Markov models

- Ergodic/Fully Connected HMMs
  - a HMM allowing for transitions from any emitting state to any other emitting state
- Left-Right HMMs
  - an HMM where the transitions are not allowed to states which indices are lower than the current state:  $a_{ij} = 0, j < i$



$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix}$$

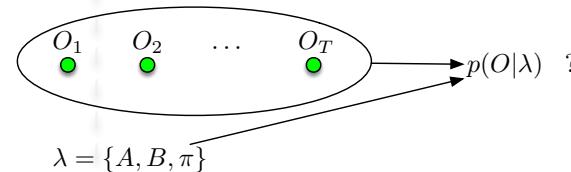


source: L. Rabiner, 1989.

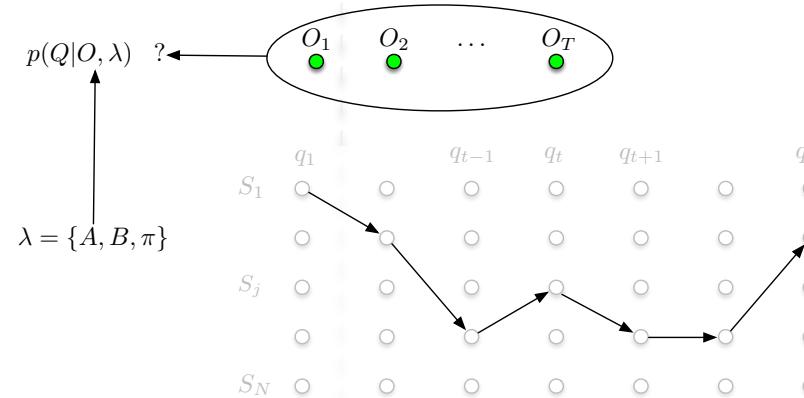
# Hidden Markov Model (HMM)

## The three basic problems for HMMs

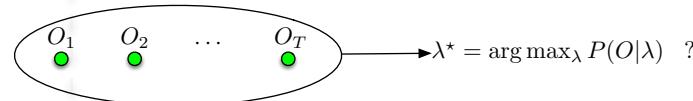
### Problem 1. Likelihood of a model (Forward algorithm)



### Problem 2. Decoding of the best sequence (Viterbi algorithm)



### Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

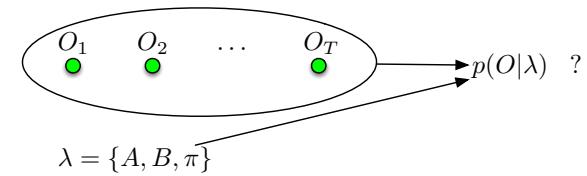


# Hidden Markov Model (HMM)

## The three basic problems for HMMs

### Problem 1. Likelihood of a model (Forward algorithm)

- Given
  - a sequence of observations  $O = O_1, O_2, \dots, O_T$
  - a model  $\lambda = \{A, B, \pi\}$
- What is the probability that the model has generated  $O$  ?
  - $P(O | \lambda)$  ?
- Application example:
  - does the observation sequence  $O$  corresponds to
    - the model  $\lambda_1$  {sleep/eat/exercice} of Doudou the hamster ?
    - the model  $\lambda_2$  {sleep/metro/work} of Bill the employee ?
- **Problem:**
  - We don't know which states  $Q = q_1, q_2, \dots, q_T$  the model has gone through
    - We need to consider all possibilities



# Hidden Markov Model (HMM)

## Problem 1. Likelihood of a model (Forward algorithm)

### Method 1 (exhaustive)

- We need to enumerate all sequences of states of length  $T$ :

$$Q = q_1, q_2, \dots, q_T$$

- For each sequence  $Q$ , we estimate

- The probability of the observations  $O$

$$P(O|Q, \lambda) = \prod_{t=1}^T P(O_t|q_t, \lambda) = b_{q_1}(O_1) \cdot b_{q_2}(O_2) \dots b_{q_T}(O_T)$$

- The probability of this specific sequence  $Q$

$$P(Q|\lambda) = \pi_{q_1} \cdot a_{q_1 q_2} \cdot a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

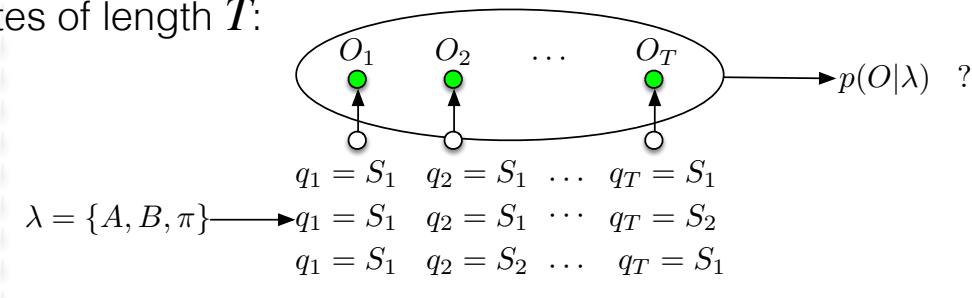
- The joint probability of  $O$  and  $Q$

$$P(O, Q|\lambda) = P(O|Q, \lambda) P(Q|\lambda) = \pi_{q_1} \cdot b_{q_1}(O_1) \cdot a_{q_1, q_2} \dots a_{q_{T-1}, q_T} \cdot b_{q_T}(O_T)$$

- Finally, we need to sum up for all possible sequences  $Q$

$$P(O|\lambda) = \sum_{\forall Q} P(O, Q|\lambda) = \sum_{\forall q_1, q_2, \dots} P(O|Q, \lambda) P(Q|\lambda)$$

- Huge computational cost !!!



# Hidden Markov Model (HMM)

## Problem 1. Likelihood of a model (Forward algorithm)

### Method 2 (optimised): Forward algorithm

#### – Definition: Forward variable

$$\alpha_t(j) \stackrel{\text{def}}{=} P(O_1, O_2, \dots, O_t, q_t = S_j | \lambda)$$

- probability of the partial observation sequence  $O_1, O_2, \dots, O_t$  (until time  $t$ ) and "state  $S_j$  at time  $t$ " given the model  $\lambda$

- Initialisation

$$\alpha_1(i) = \pi_i \cdot b_i(O_1)$$

- Loop

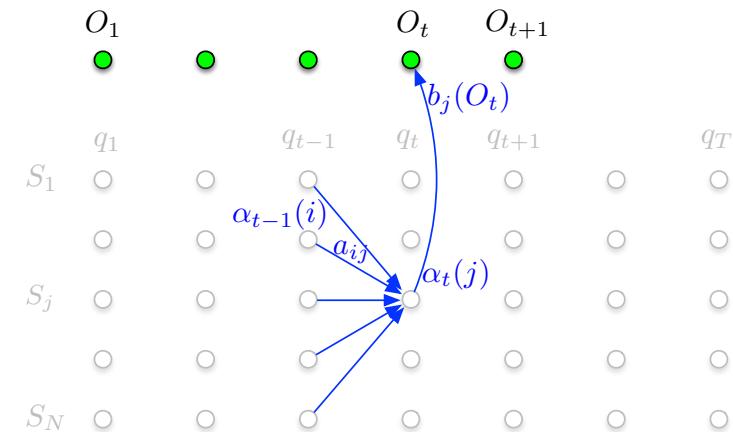
$$\alpha_t(j) = \left[ \sum_{i=1}^N \alpha_{t-1}(i) \cdot a_{ij} \right] b_j(O_t)$$

- Ending

$$P(O | \lambda) = \sum_{i=1}^N \alpha_T(i)$$

Proof:

$$\begin{aligned}\alpha_T(i) &= P(O_1, O_2, \dots, O_T, q_t = S_i | \lambda) \\ \sum_i \alpha_T(i) &= \sum_i P(O_1, O_2, \dots, O_T, q_t = S_i | \lambda) \\ &= P(O_1, O_2, \dots, O_T | \lambda)\end{aligned}$$



# Hidden Markov Model (HMM)

## Problem 1. Likelihood of a model (Forward algorithm)

### Method 2 (optimised): revisited

- We could also have used the reverse recursion

- **Definition: Backward variable**

$$\beta_t(i) \stackrel{\text{def}}{=} P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = S_i, \lambda)$$

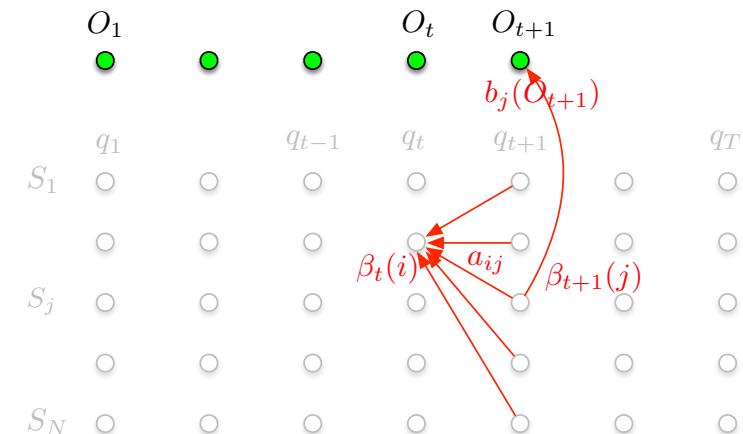
- probability of the partial observation sequence from  $t + 1$  to  $T$  given "state  $S_i$  at time  $t$ ", and the model  $\lambda$

- Initialisation

$$\beta_T(i) = 1$$

- Loop

$$\beta_t(i) = \sum_j a_{ij} \cdot b_j(O_{t+1}) \cdot \beta_{t+1}(j)$$



# Hidden Markov Model (HMM)

## The three basic problems for HMMs

### Problem 2. Decoding of the best sequence (Viterbi algorithm)

– Given

- a sequence of observations  $O = O_1, O_2, \dots, O_T$
- a model  $\lambda = \{A, B, \pi\}$

– What is the corresponding sequence of states  $Q = q_1, q_2, \dots, q_T$  ?

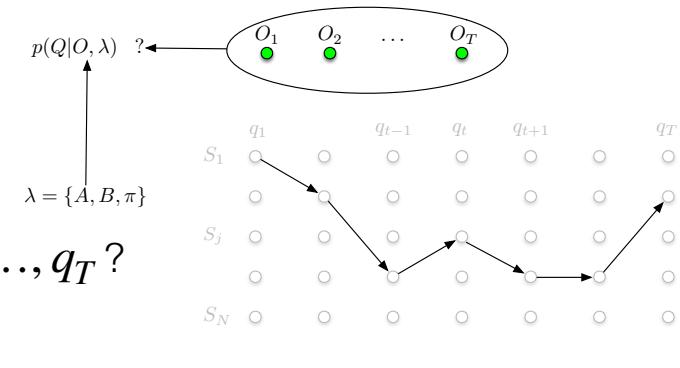
- $Q^* = \arg \max_Q P(Q | O, \lambda)$  ?

– Application example:

- if we observe the sequence of sounds  $O$  produced by Doudou and given its model  $\lambda$  {sleep/eat/exercice}, what is the corresponding sequence of activities  $Q$  of Doudou ?

– **Problem:**

- Find the most **optimal** sequence  $Q$ 
  - Several possible definitions for **optimality**



# Hidden Markov Model (HMM)

## Problem 2. Decoding of the best sequence (Viterbi algorithm)

**Optimality-1: choose the  $q_t$  which are individually the most likely**

– Définition

$$\gamma_t(i) \stackrel{\text{def}}{=} P(q_t = S_i | O, \lambda)$$

- probability of "being in state  $S_i$  at time  $t$ " given the observation sequence  $O$  and the model  $\lambda$
- using  $p(q | O, \lambda) p(O | \lambda) = p(q, O | \lambda)$

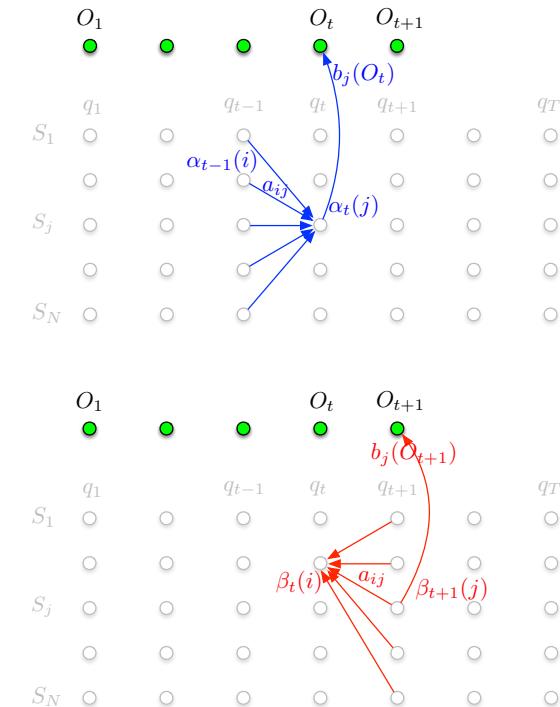
Reminder:

$$\begin{aligned}\alpha_t(i) &= P(O_1, O_2, \dots, O_t, q_t = S_i | \lambda) \\ \beta_t(i) &= P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = S_i, \lambda)\end{aligned}$$

$$\gamma_t(i) = \frac{p(q_t = S_i, O | \lambda)}{p(O | \lambda)} = \frac{\alpha_t(i)\beta_t(i)}{\sum_q p(q, O | \lambda)} = \frac{\alpha_t(i)\beta_t(i)}{\sum_i \alpha_t(i)\beta_t(i)}$$

– We choose  $q_t = \arg \max_i [\gamma_t(i)]$

- Problem: it does not take into account the transition probabilities (what about if  $a_{ij} = 0$ ?)



# Hidden Markov Model (HMM)

## Problem 2. Decoding of the best sequence (Viterbi algorithm)

**Optimality-2: choose the  $Q$  which is globally the most likely**

### – Dynamic Programming

- $\delta_t(j) \stackrel{\text{def}}{=} \max_{q_1, q_2, \dots, q_{t-1}} P(O_1, O_2 \dots O_t, q_1, q_2 \dots q_t = S_j | \lambda)$

- it is the highest probability of the partial alignment starting at 1 and ending at  $t$  in state  $S_j$

- We can compute  $\delta_t(*)$  using  $\delta_{t-1}(*)$  as

$$\delta_t(j) = \max_i [\delta_{t-1}(i) a_{ij}] b_j(O_t)$$

### – Viterbi algorithm

- **Forward pass**

- Initialisation:

$$\delta_1(i) = \pi_i b_i(O_1)$$

- Loop

for  $t=2:T$  and for  $j=1:N$

$$\delta_t(j) = \max_i (\delta_{t-1}(i) a_{ij}) b_j(O_t)$$

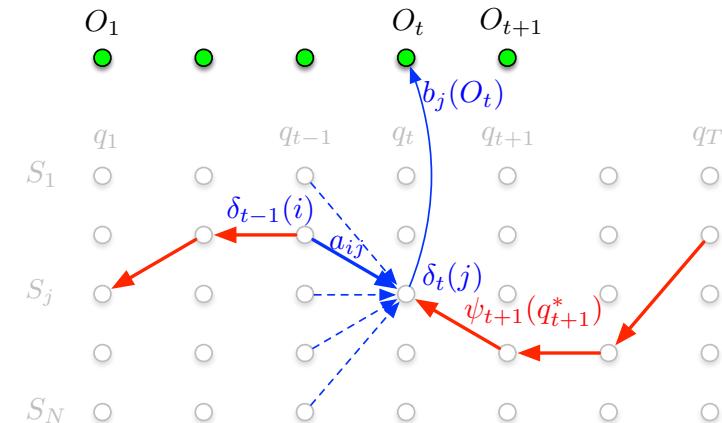
- Ending

$$P^* = \max_i \delta_T(i)$$

- **Backward pass**

- for  $T=T-1, T-2, \dots, 1$

$$q_t^* = \psi_{t+1}(q_{t+1}^*)$$



and  $\psi_1(i) = 0$

and  $\psi_t(j) = \arg \max_i (\delta_{t-1}(i) a_{ij})$

and  $q_T^* = \arg \max_i \delta_T(i)$

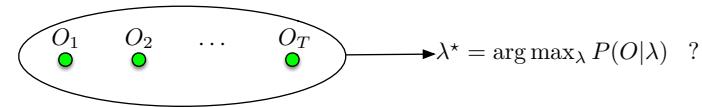


# Hidden Markov Model (HMM)

## The three basic problems for HMMs

### Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

- Given
  - a sequence of observations  $O = O_1, O_2, \dots, O_T$
- What is the model  $\lambda^*$  that maximises the likelihood of the observations:
  - $\lambda^* = \arg \max_{\lambda} P(O | \lambda)$
  - i.e. find the parameters  $\lambda = \{A, B, \pi\}$  that maximises the likelihood of the observation's sequence  $O$
- Application example:
  - find the parameters of the model  $\lambda$  of Doudou the hamster ?



## Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

### Baum-Welch algorithm

- EM (Expectation-Maximisation) algorithm

- Definition

$$\xi(i, j) \stackrel{\text{def}}{=} P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

- using  $p(q | O, \lambda) p(O | \lambda) = p(q, O | \lambda)$

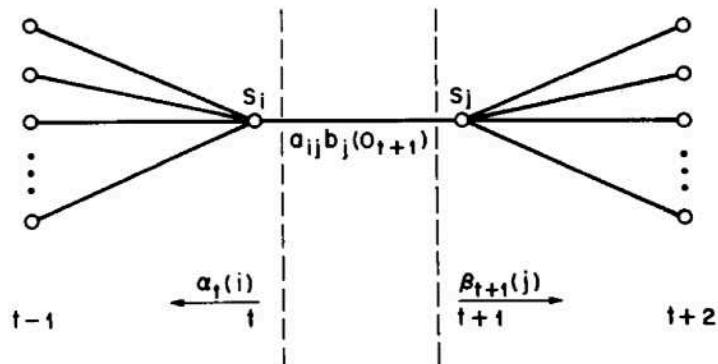
$$\begin{aligned}\xi(i, j) &= \frac{P(q_t = S_i, q_{t+1} = S_j, O | \lambda)}{P(O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_i \sum_j \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}\end{aligned}$$

- We also have

$$\gamma_t(i) = \sum_j \xi_t(i, j)$$

Reminder:

$$\begin{aligned}\alpha_t(i) &= P(O_1, O_2, \dots, O_t, q_t = S_i | \lambda) \\ \beta_t(i) &= P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = S_i, \lambda) \\ \gamma_t(i) &= P(q_t = S_i | O, \lambda)\end{aligned}$$



**Fig. 6.** Illustration of the sequence of operations required for the computation of the joint event that the system is in state  $S_i$  at time  $t$  and state  $S_j$  at time  $t + 1$ .

source: L. Rabiner, 1989.

## Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

### Baum-Welch algorithm (cont.)

- $\sum_{t=1}^{T-1} \gamma_t(i) =$  expected (over time) number of times that  $S_i$  is visited
- $\sum_{t=1}^{T-1} \xi_t(i, j) =$  expected (over time) number of transitions from  $S_i$  to  $S_j$

#### – Re-estimation of the HMM parameters

- $\pi_i =$  expected frequency (number of times) in  $S_i$  at time ( $t = 1$ ) =  $\gamma_1(i)$

- $a_{ij} = \frac{\text{expected number of transitions from state to state } S_i \text{ to state } S_j}{\text{expected number of transitions from state } S_i} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$

#### – a) Discrete observations

- $b_j(k) = \frac{\text{expected number of times in state } S_j \text{ an observing symbol } v_k}{\text{expected number of times in state } S_j} = \frac{\sum_{t=1}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$

## Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

### Baum-Welch algorithm (cont.)

#### – a) Discrete observations

$$\bullet \quad b_j(k) = \frac{\text{expected number of times in state } S_j \text{ an observing symbol } v_k}{\text{expected number of times in state } S_j} = \frac{\sum_{t=1 s.t. O_t=v_k}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$$



## Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

### Baum-Welch algorithm (cont.)

#### – b) Continuous Observation Densities in HMMs

- Gaussian Mixture Model with  $m \in \{1, \dots, M\}$  components for state  $S_j$

$$b_j(O) = \sum_{m=1}^M c_{jm} \cdot \mathcal{N}(O, \mu_{jm}, \Sigma_{jm})$$

- $\gamma_t(j, k)$  is the probability of being in state  $j$  at time  $t$  with the  $k^{th}$  mixture component

$$\gamma_t(j, k) = \left[ \frac{\alpha_t(j)\beta_t(j)}{\sum_{j=1}^N \alpha_t(j)\beta_t(j)} \right] \cdot \left[ \frac{c_{jk}\mathcal{N}(O_t, \mu_{jk}, \Sigma_{jk})}{\sum_{m=1}^M c_{jm}\mathcal{N}(O_t, \mu_{jm}, \Sigma_{jm})} \right]$$

- Re-estimating HMM parameters for state  $j$  and  $k^{th}$  mixture component

$$\text{– mixture coefficients: } c_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k)}{\sum_{t=1}^T \sum_{k=1}^M \gamma_t(j, k)}$$

$$\text{– mean vectors: } \mu_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k) \cdot O_t}{\sum_{t=1}^T \gamma_t(j, k)}$$

$$\text{– covariance matrice: } \Sigma_{jk} = \frac{\sum_{t=1}^T \gamma_t(j, k) \cdot (O_t - \mu_{jm})(O_t - \mu_{jm})'}{\sum_{t=1}^T \gamma_t(j, k)}$$

## Problem 3. Training a HMM model (Forward-Backward, Baum-Welch algorithm)

### Baum-Welch algorithm (cont.)

#### – HMM parameters constraints

$$\sum_{i=1}^N \pi_i = 1$$

$$\sum_{j=1}^N a_{ij} = 1, \quad 1 \leq i \leq N$$

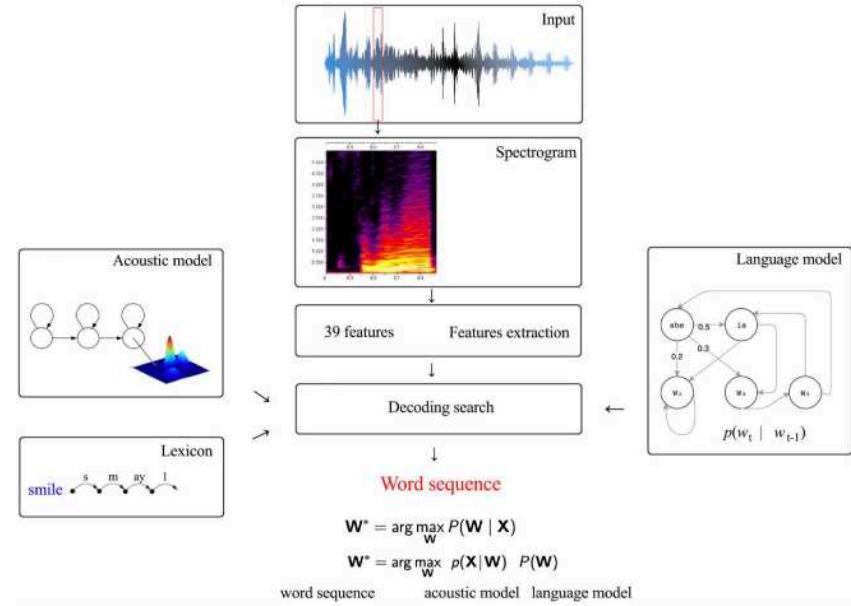
$$\sum_{k=1}^M b_j(k) = 1, \quad 1 \leq j \leq N$$

# Automatic Speech Recognition (ASR) using GMM/HMM

# Automatic Speech Recognition (ASR) using GMM/HMM

## Automatic Speech Recognition using GMM/HMM

- Given a speech signal  $s$
- Find the most likely word sequence
  - $p(W|X) = \frac{p(X|W) \cdot p(W)}{p(X)}$
- $p(W)$ :
  - prior probability of the word sequence  $w$  obtained from a **Language model**
- $p(X|W)$ 
  - likelihood of the word sequence  $w$  obtained from stored models of speech (**Acoustic model**)



source: <https://jonathan-hui.medium.com/speech-recognition-gmm-hmm-8bb5eff8b196>

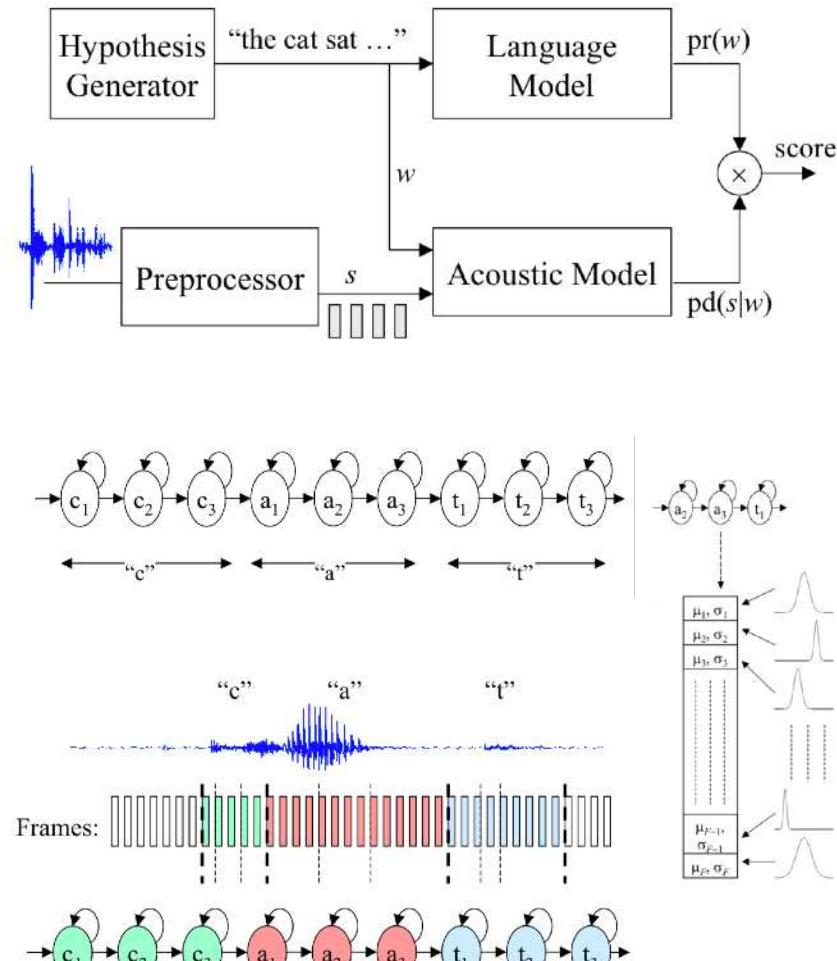
# Automatic Speech Recognition (ASR) using GMM/HMM

## Automatic Speech Recognition using GMM/HMM

### – Isolated Word Recognition

- each phoneme in a word corresponds to a number of model states (typically 3 states per phoneme)
- when saying a word, the speaker stays in each state for one or more frames and then goes on to the next state
- the emission probability of each state is modeled as a GMM over the MFCC/ $\Delta/\Delta^2$
- $p(s | w)$  is obtained by Viterbi alignment
- choose the word with the highest probability
- separate HMM for each word in the vocabulary

### – Continuous Speech Recognition



source: Brooke, 2002, Speech Processing





# Lecture What you need to know

- What is a Source/Filter model
- What are MFCCs ?
- What are the Chroma/ Pitch-Class-Profile
  - how they are computed
  - what are there advantages/ drawbacks
- What is the Constant-Q-Transform
  - what problem does it solve?
  - how does it solve it ?
- What is a pitch
  - what is a temperament (tuning-system)
- What is a Hidden Markov Model
  - what are the variables
  - what is the Viterbi decoding algorithm
- How are chords estimated using an HMM approach, a DL approach ?
- How are the performances of such a system measured ?