

M2 Data Science

DS-telecom-15 "Audio and Music Information Retrieval"



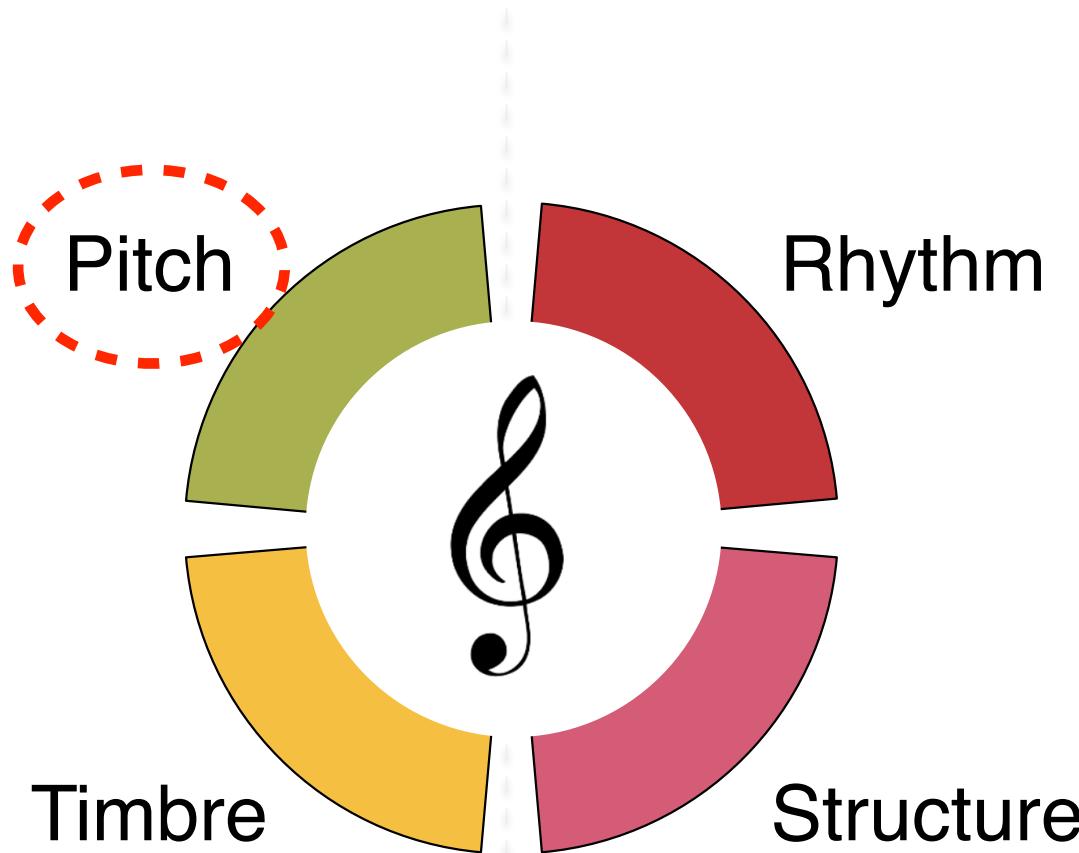
Geoffroy Peeters

contact: geoffroy.peeters@telecom-paris.fr

Télécom-Paris, IP-Paris, France

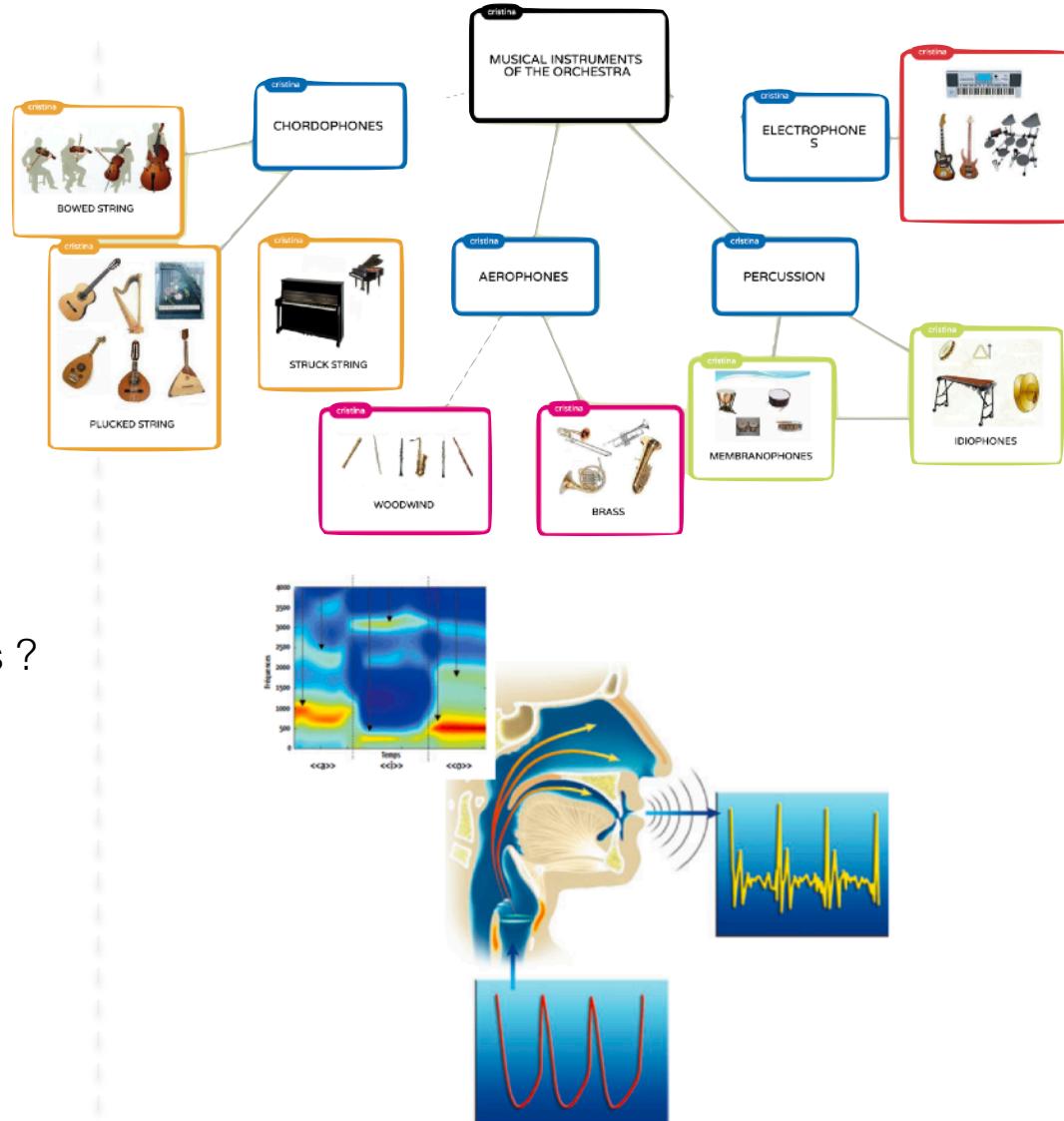
Reminder of Musical Concepts

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

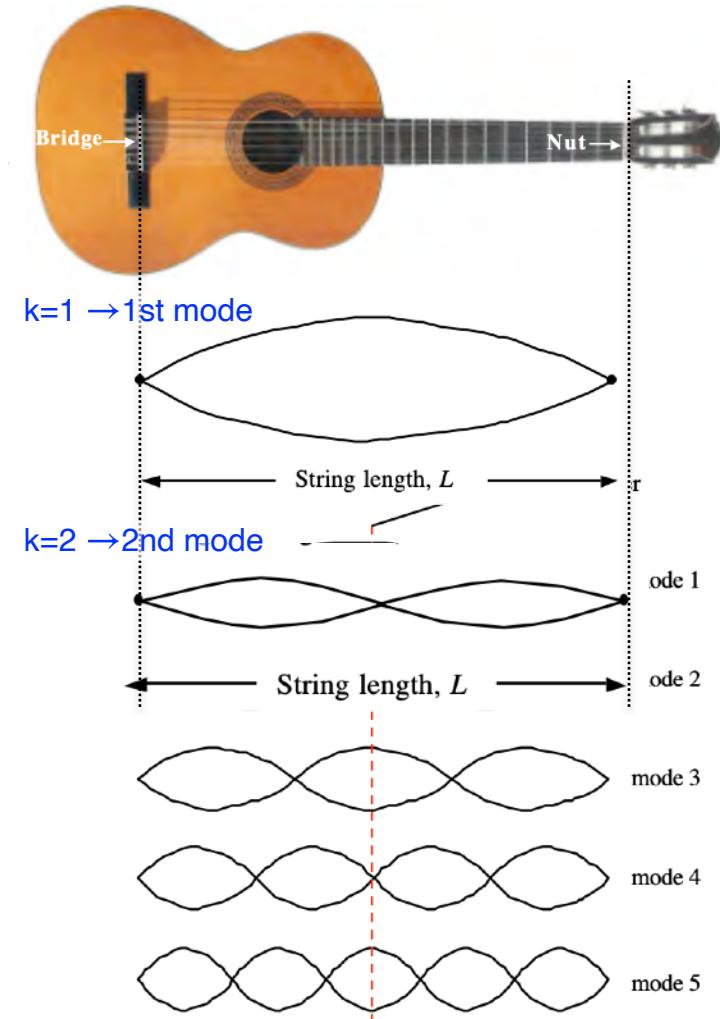
Harmonic sounds

- **Reminder**

- Frequency of a fingered string

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

- T : string tension,
- μ : linear mass density,
- k : vibration mode or harmonic number



Harmonic sounds

- Harmonic representation of sound
 - Fourier series

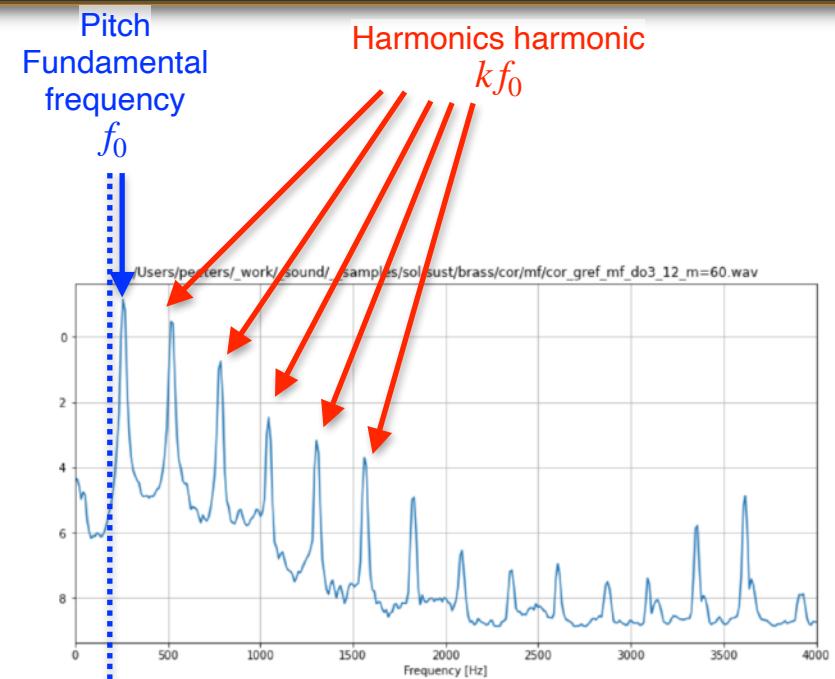
$$f(t) = \sum_{n=0}^{+\infty} k_n \sin(2\pi n f_0 t + \phi_n)$$

$$= \sum_{n=0}^{+\infty} A_n \cos(2\pi n f_0 t) + B_n \sin(2\pi n f_0 t)$$

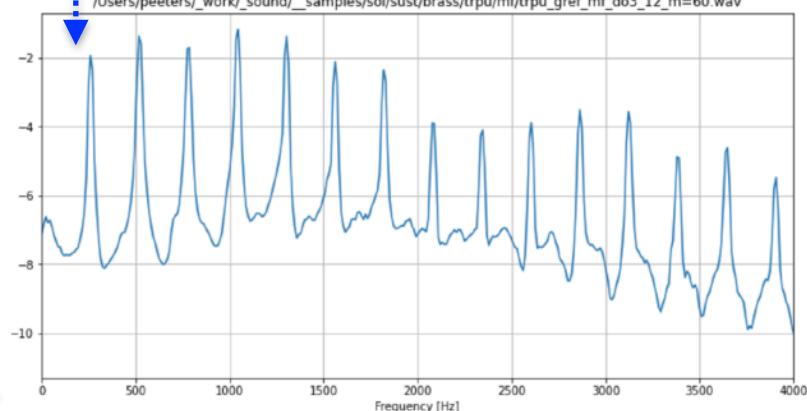
$$= \sum_{n=-\infty}^{+\infty} c_n e^{j2\pi n f_0 t}$$

- 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



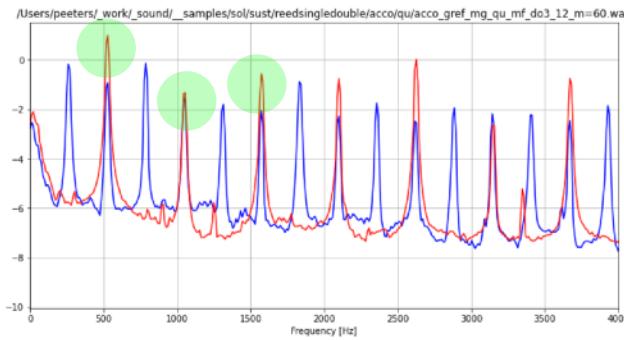
Same pitch
Values (amplitude) of the harmonics define the "timbre"



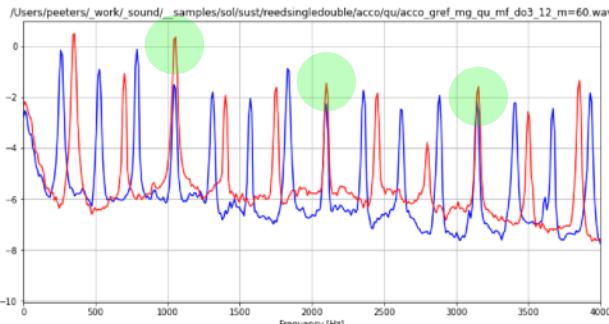
Harmonic sounds

- Consonance/ Dissonance between pitches
 - can be explained by the overlap between the harmonics
 - closeness between harmonics that will fall into a critical band

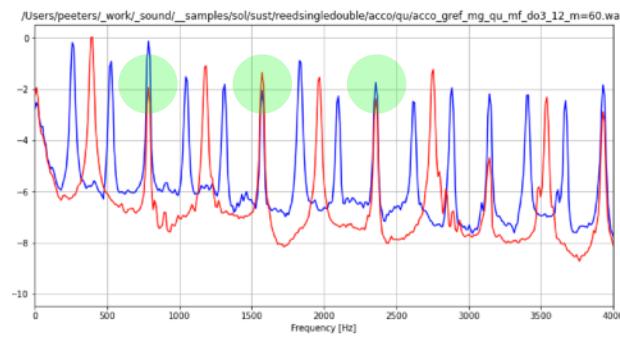
Octave interval ($F_0 = 2f_0$)



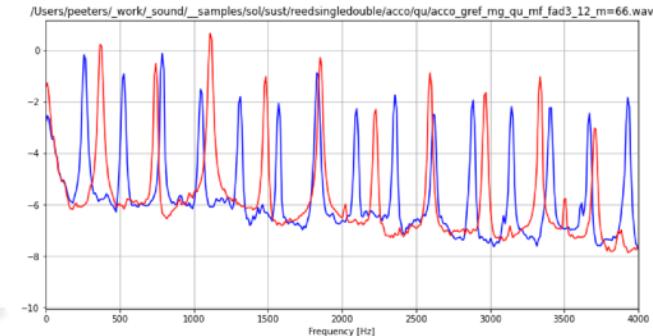
4th interval ($3F_0 = 4f_0$)



5th interval ($2F_0 = 3f_0$)



Triton (most dissonant !!!)

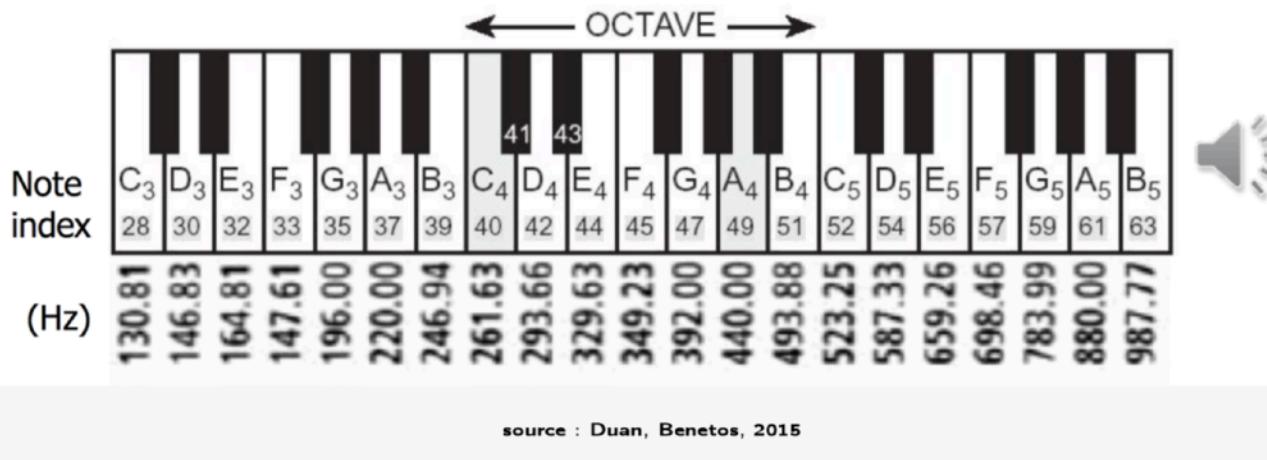


What is pitch ?

- **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, ré-3, mi-3



What is pitch ?

Consonance

- To map Pitch to Notes we need to define an organization of the pitches
 - The **temperament** is a **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 (the most consonant) ... because of the harmonic series -> 2nd harmonics
 - **Fifth**: do-3 → sol-3 (very consonant) ... because of the harmonic series -> 3rd harmonics
 - **Fourth**: do-3 → fa-3 (very consonant) ... because of the harmonic series -> 4th 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

$$f = \alpha \frac{1}{2L}$$

- If string fingered

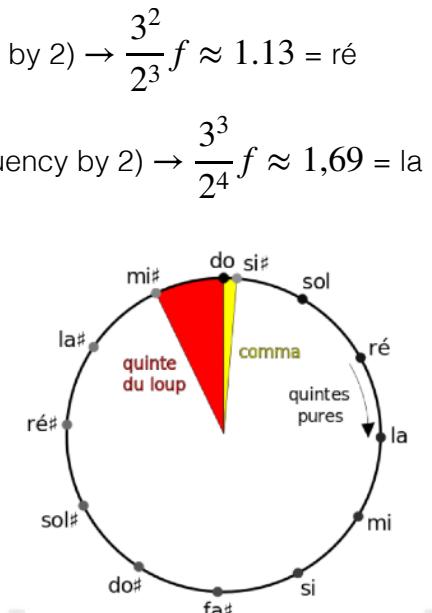
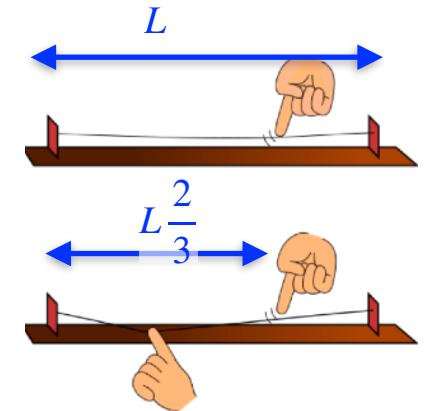
- at $L \left(\frac{1}{2}\right) \rightarrow 2f$ → octave (consonant interval)
- at $L \left(\frac{2}{3}\right) \rightarrow \frac{3}{2}f = 1.5f$ → 5th (consonant interval) = sol
- at $L \left(\frac{2}{3}\right)^2 \rightarrow \frac{3^2}{2^2}f$ → 5th (5th) → bring it back to main octave (divide frequency by 2) → $\frac{3^2}{2^3}f \approx 1.13$ = ré
- at $L \left(\frac{2}{3}\right)^3 \rightarrow \frac{3^3}{2^3}f$ → 5th (5th (5th)) → bring it back to main octave (divide frequency by 2) → $\frac{3^3}{2^4}f \approx 1.69$ = la
- and so one

- Resulting scale

- 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 !!!



What is pitch ? Temperaments

- **(2) Equal temperament**

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

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

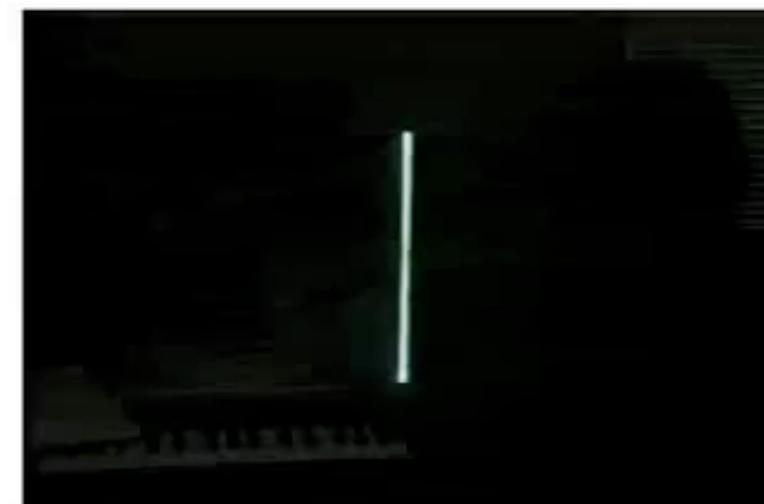
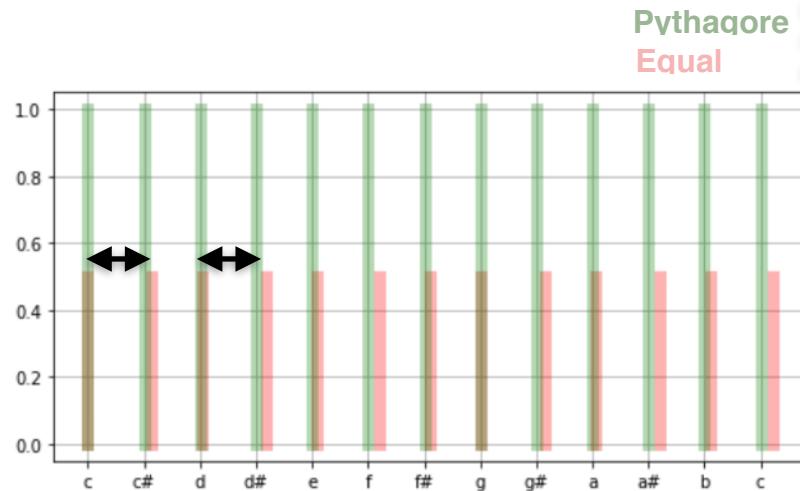
$$f_1 = f \cdot 2^{\frac{1}{12}}$$

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

...

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

- **Comparing Pythagorean and Equal Temperament**



https://fr.wikipedia.org/wiki/Courbe_de_Lissajous

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

What is pitch ?

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

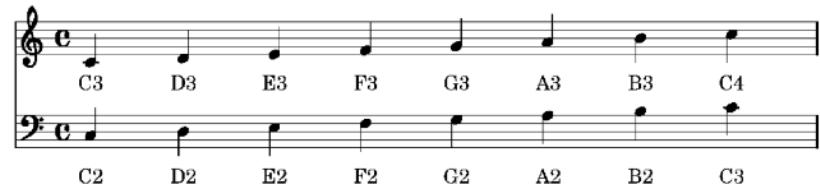
- 261 Hz → C-4 or do-3 (do at 3rd octave)
- 440 Hz → A-4 or la-3 (la at 3rd octave)
- 880 Hz → A-5 or la-4 (la at 4th octave)

- midi values

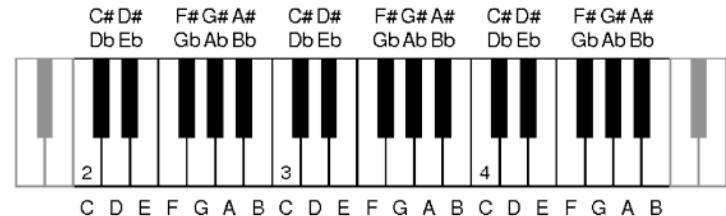
- $\frac{f_1}{f_0} = 2^{\frac{1}{12}} \rightarrow f = 440 \cdot 2^{\frac{(m-69)}{12}}$
- $m = 12 \log_2 \left(\frac{f}{440} \right) + 69$
- 261 Hz → 60
- 440 Hz → 69
- 880 Hz → 81

Note	C3	D3	E3	F3	G3	A3	B3	C4
f_0 , MIDI	48	50	52	53	55	57	59	60
f_0 , Hz	130.8	146.8	164.8	174.6	196.0	220.0	246.9	261.6

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



(b) Note pitches on the common musical notation.



What is pitch ?

Pitch representation

- Representation of pitch: musical score



Time

key/tonality

tempo

time signature

Piano

Andantino

con Pedale

Instrument (left/right hand)

bars separation

intensity

rit. e dim..... pp

PRELUDE

Op. 26, No. 7

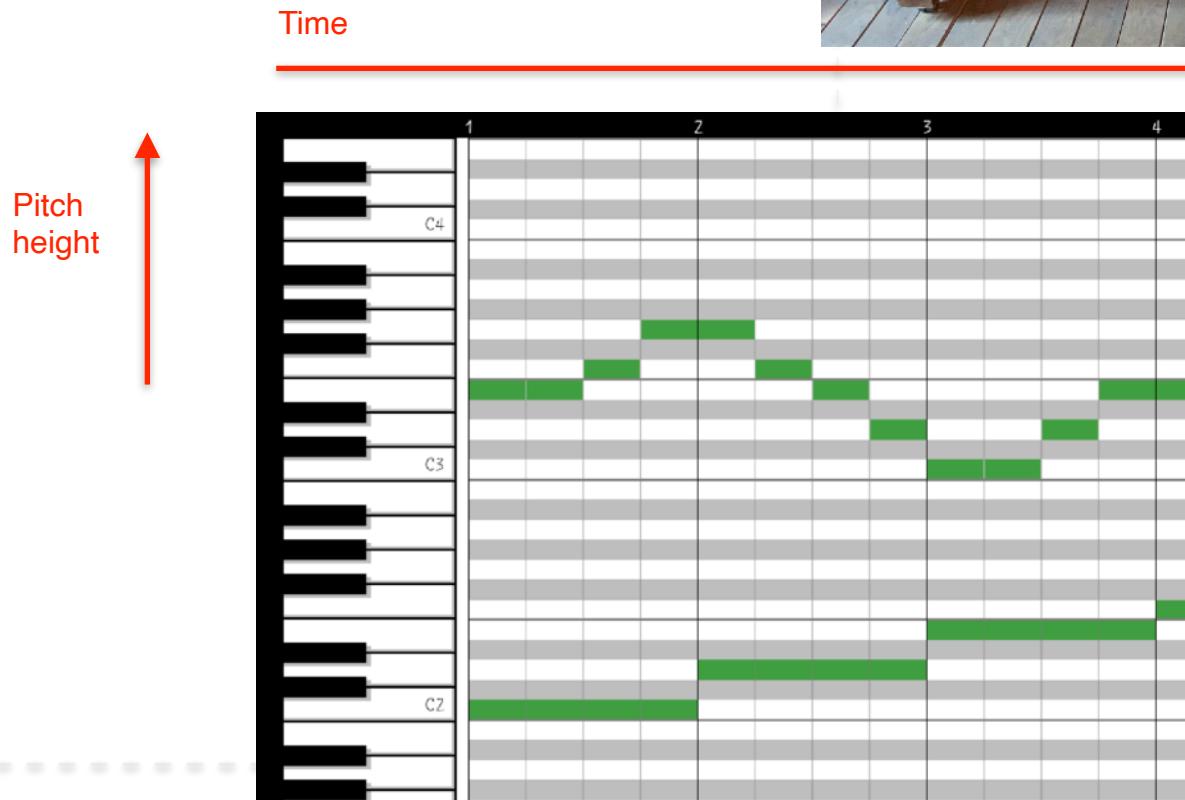
Frederic Chopin

This diagram illustrates the various elements of a musical score, specifically a piano prelude by Frederic Chopin. It features two staves of musical notation. A red arrow labeled 'Time' points horizontally across the top. A vertical red arrow labeled 'Pitch height' points upwards from the bottom staff. A green dashed circle labeled 'key/tonality' encloses the treble clef and key signature. A blue dashed circle labeled 'tempo' encloses the instruction 'Andantino'. A blue dashed circle labeled 'time signature' encloses '3/4'. A bracket labeled 'Instrument (left/right hand)' spans both staves. Blue arrows point to specific markings: 'con Pedale' on the first staff, 'bars separation' between the staves, and 'intensity' on the second staff. The score is titled 'PRELUDE Op. 26, No. 7' and is attributed to 'Frederic Chopin'.

What is pitch ?

Pitch representation

- **Representation of pitch: piano roll**



What is pitch ?

How pitches are organized ?

- **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



- **Homophony**

- 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



What is pitch ?

How pitches are organized ?

- **Modal music**

- since the Greeks, Gregorian and Church
- which uses modal scales
 - (as opposed to Major/minor scales)
- centered 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 two sets of musical notation on a single staff. The top set, enclosed in a red box, is labeled "Ionian Major". It consists of six notes: A, B, C, D, E, F, G. The bottom set, also enclosed in a red box, is labeled "Aeolian Minor (natural)". It consists of seven notes: A, B, C, D, E, F, G. Both sets of notes are separated by vertical bar lines.

Greek, Gregorian, Church modes

The image displays a series of musical staves, each representing a different mode or raga. The modes listed are: Bilâwal, Kalyan, Bhairav, Purvi, Mârvâ, Bhairavi, Åsâvâri, Kâfi, and Todi. Each mode is associated with a specific sequence of notes (Sa, Re, Ga, Ma, Pa, Dha, Ni, Sa) and a unique pattern of note heads (solid, open, or filled circles). The notes are separated by vertical bar lines.

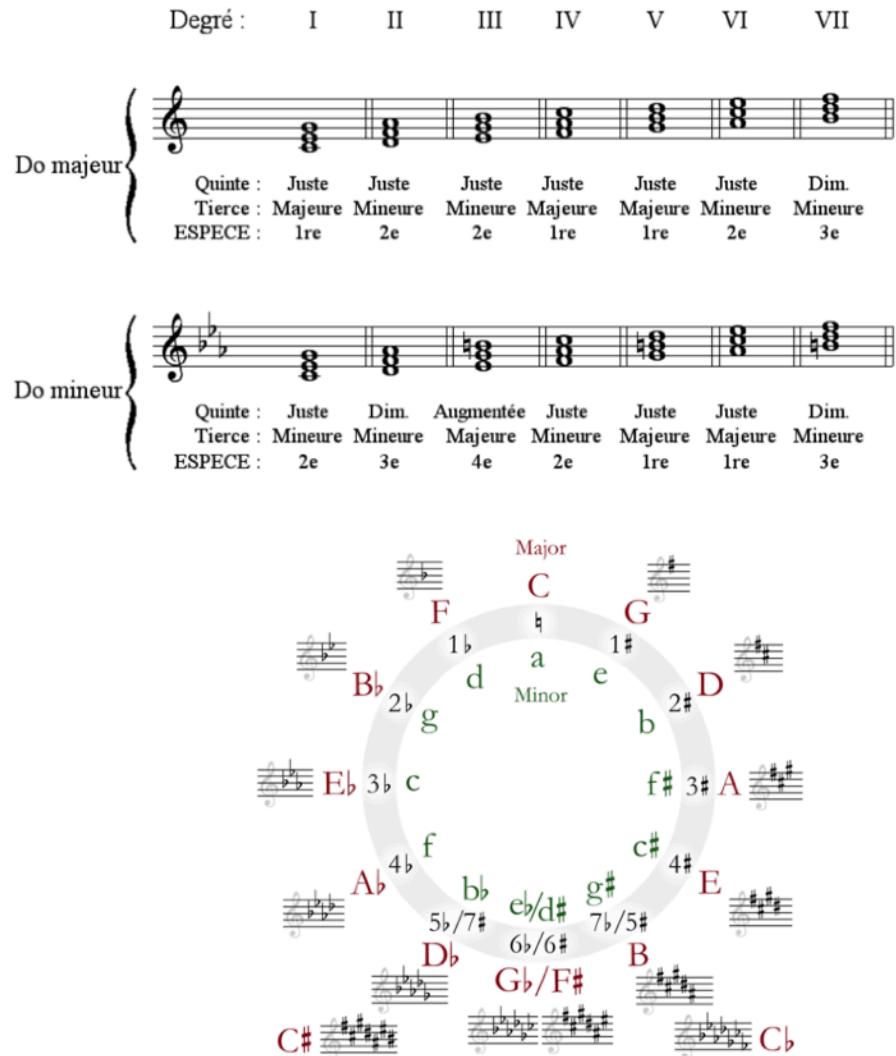
Indian raga

What is pitch ?

How pitches are organized ?

- **Tonal music**

- centered 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 of often, notes are played (almost) simultaneously
 - but also can be arpeggiated

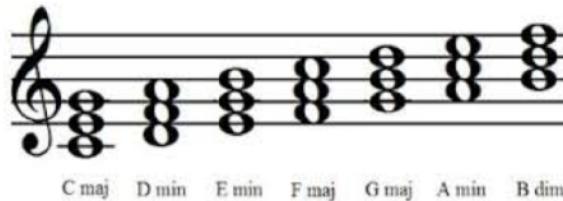
- Chords on Major scale

I II III IV V VI VII VIII

C MAJ⁷ D MIN⁷ E MIN⁷ F MAJ⁷ G⁷ A MIN⁷ B MIN^{7(b5)} C MAJ⁷

Chord dictionary

Root	Major	m	+	6	m6	7	m7	ma7	o	9
C	C	Cm	C+	C6	Cm6	C7	Cm7	Cma7	C°	C9
C#	D	D#	C#m	D#+	D#6	C#m6	D#7	C#ma7	D#P	D#9
D	D	D	Dm	D+	D6	Dm6	D7	Dma7	D°	D9
D#	E	E#	E#m	E#+	E#6	E#m6	E#7	E#ma7	E#P	E#9
E	F	F#	F#m	F#+	F#6	F#m6	F#7	F#ma7	F#P	F#9
F	G	G#	G#m	G#+	G#6	G#m6	G#7	G#ma7	G°	G9
F#	G#	G#	G#m	G#+	G#6	G#m6	G#7	G#ma7	G#P	G#9
G	A	A#	A#m	A#+	A#6	A#m6	A#7	A#ma7	A°	A9
A	B	B#	B#m	B#+	B#6	B#m6	B#7	B#ma7	B°	B9
B	C	C#	C#m	C#+	C#6	C#m6	C#7	C#ma7	C°	C9
Root-note	Type of the chord									
c, c#, d, d# ... b	Triads: major (C-M: c,e,g) , minor (C-m: c, e♭, g), suspended (C-sus2: c, d, g / C-sus4: c, f, g), augmented (C-aug: c, e, g#), diminished (C-dim: c, e♭, g♭)									
	Tetrads: major 7 (C-M7: c, e, g, b), minor 7 (C-m7: c, e♭, g, b♭), dominant 7 (C-7: c, e, g, b), major 6 (C-M6: c, e, g, a), minor 6 (C-m6: c, e♭, g, a) ...									
	Pentads: major 9 (C-M9: c, e, g, b, d), dominant 9 (C-9: c, e, g, b♭, d) ...									



C maj D min E min F maj G maj A min B dim

Automatic Chord Recognition (ACR)

Automatic Chord Recognition (ACR)

audio / music



Automatic Chord Recognition System

sequence of chords

A A AAA AF#
I need another story

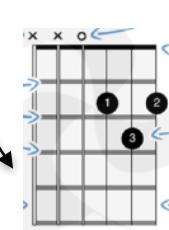
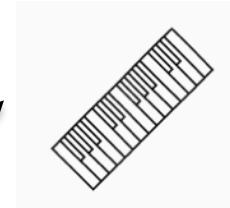
AA AA AE E
Something to get off my chest

A A A AA AF#
My life gets kinda boring

A AA AA A A D DB
Need something that I can confess

A A B D E C# D
'Til all my sleeves are stained red

A A B C# D C# C#
From all the truth that I've said



guitar tab

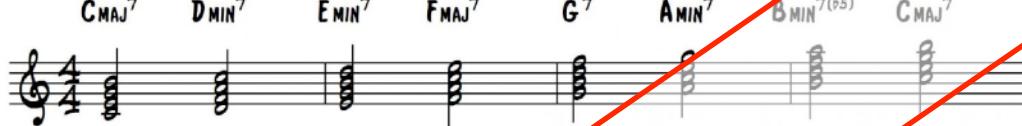
What is pitch ?

→ How pitches are organized ?

- **How pitches are organized ?**

→ **les accords:**

- ensemble de notes considéré comme formant un tout du point de vue de l'harmonie
- le plus souvent, les notes sont simultanées (superposées en un même moment)
 - mais aussi, notes successives, par exemple dans des arpèges.



- Accords sur les degrés de la gamme majeur

Dictionnaire d'accords

Root	Major	m	+	6	m6	7	m7	ma7	o	9
C	C	Cm	C+	C6	Cm6	C7	Cm7	Cma7	C°	C9
C#	D	D#	C#m	D#+	D#6	C#m6	D#7	C#m7	D#ma7	D#9
D	D	D	Dm	D+	D6	Dm6	D7	Dm7	Dma7	D9
D#	E	E#	E#m	E#+	E6	E#m6	E7	E#m7	E#ma7	E#9
E	F	F#	F#m	F#+	F6	F#m6	F7	F#m7	F#ma7	F#9
F	G	G#	G#m	G#+	G6	G#m6	G7	G#m7	G#ma7	G#9
F#	G#	G	Gm	G+	G6	Gm6	G7	Gm7	Gma7	G9
G	A	A#	A#m	A#+	A6	A#m6	A7	A#m7	A#ma7	A#9
A	B	B#	B#m	B#+	B6	B#m6	B7	B#m7	B#ma7	B#9
B	C	C#	C#m	C#+	C6	C#m6	B7	B#m7	B#ma7	B#9
C	D	D#	D#m	D#+	D6	D#m6	C7	C#m7	C#ma7	C#9
D#	E	F	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7	E#9
E	F#	G	G#m	G#+	G6	G#m6	F7	F#m7	F#ma7	F#9
F	G#	A	A#m	A#+	A6	A#m6	E7	E#m7	E#ma7	E#9
G#	A#	B	B#m	B#+	B6	B#m6	D7	D#m7	D#ma7	D#9
A#	B#	C	C#m	C#+	C6	C#m6	C7	C#m7	C#ma7	C#9
B#	C#	D	D#m	D#+	D6	D#m6	B7	B#m7	B#ma7	B#9
C#	D#	E	E#m	E#+	E6	E#m6	E7	E#m7	E#ma7	E#9
D	E#	F#	G	G#m	G#+	G6	G#m6	G7	G#m7	G#ma7
E	F	G#	A#	A#m	A#+	A6	A#m6	F7	F#m7	F#ma7
F	G	A#	B#	B#m	B#+	B6	B#m6	E7	E#m7	E#ma7
G	A	B#	C#	C#m	C#+	C6	C#m6	D7	D#m7	D#ma7
A	B	C#	D#	D#m	D#+	D6	D#m6	B7	B#m7	B#ma7
B	C	D#	E#	E#m	E#+	E6	E#m6	E7	E#m7	E#ma7
C	D	E#	F#	F#m	F#+	F6	F#m6	C7	C#m7	C#ma7
D	E	F#	G#	G#m	G#+	G6	G#m6	E7	E#m7	E#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#	D#m	D#+	D6	D#m6	D7	D#m7	D#ma7
B#	C#	D#	E#	E#m	E#+	E6	E#m6	B7	B#m7	B#ma7
C#	D#	E#	F#	F#m	F#+	F6	F#m6	E7	E#m7	E#ma7
D	E#	F#	G#	G#m	G#+	G6	G#m6	C7	C#m7	C#ma7
E	F#	G#	A#	A#m	A#+	A6	A#m6	D7	D#m7	D#ma7
F#	G#	A#	B#	B#m	B#+	B6	B#m6	C7	C#m7	C#ma7
G#	A#	B#	C#	C#m	C#+	C6	C#m6	E7	E#m7	E#ma7
A#	B#	C#	D#							

Automatic Chord Recognition (ACR)

Dictionnaire d'accords, vocabulaires d'accords

– 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)

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

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

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

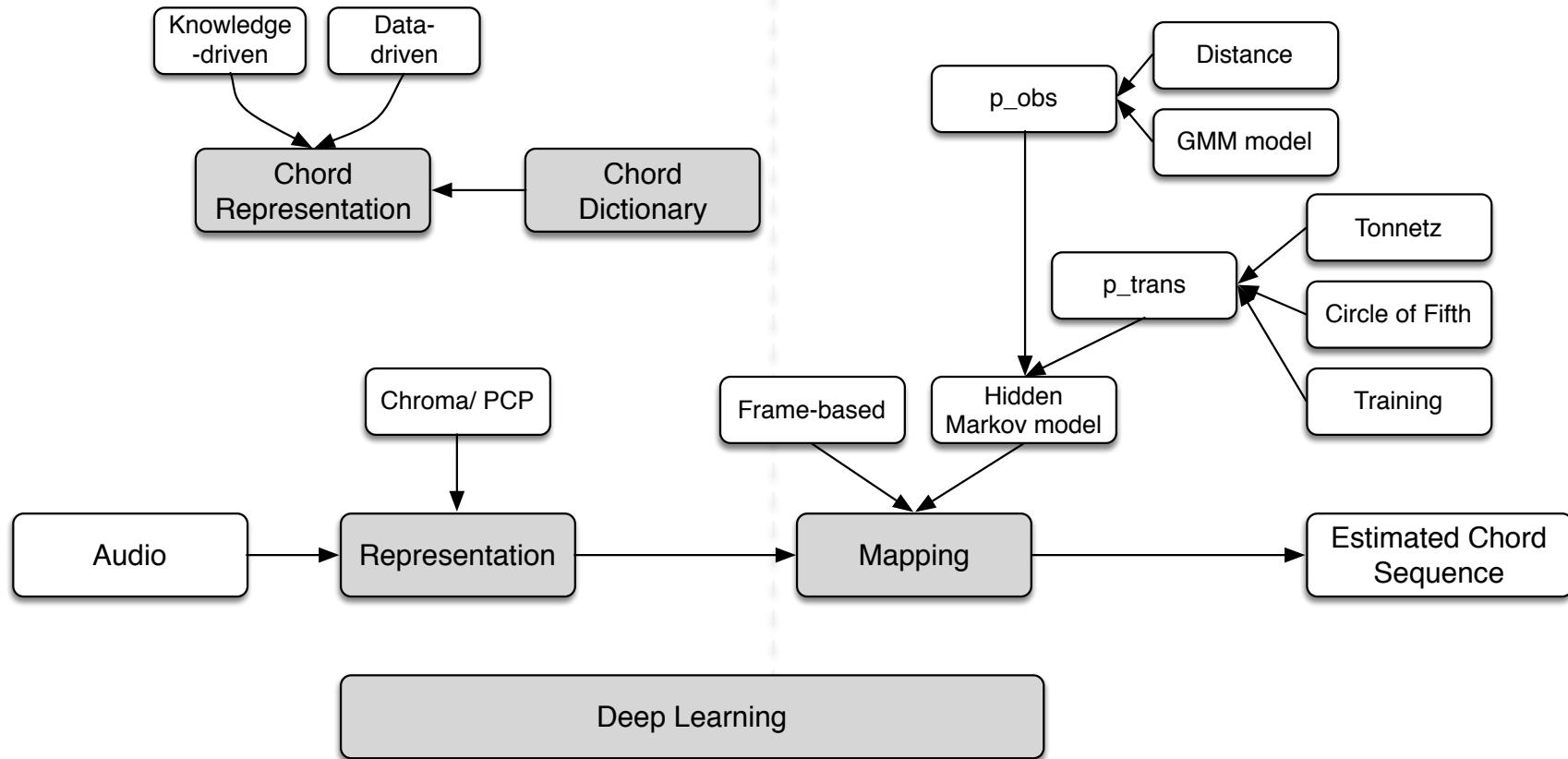
$$A/3 \equiv A:maj/3 \equiv A:(3,5)/3$$

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

Automatic Chord Recognition (ACR) Systems

Automatic Chord Recognition (ACR) Systems

Automatic Chord Recognition System



Automatic Chord Recognition (ACR) Systems

Brief overview of system evolution

- Frame-based/ template-based approach
 - **1999** → 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.
- 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 (International Society for Music Information Retrieval), pages 183–189, Baltimore, Maryland, USA, 2003
 - **2007** → H. Papadopoulos and G. Peeters. Large-scale study of chord estimation algorithms based on chroma representation. In Proc. of IEEE CBMI (International Workshop on Content-Based Multimedia Indexing), Bordeaux, France, 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 Transactions on Audio, Speech, and Language Processing*, 20(6):1771–1783, 2012.
- Deep learning approaches
 - **2013** → Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent. Audio chord recognition with recurrent neural networks. In ISMIR, pages 335–340, 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 (International Society for Music Information Retrieval), Suzhou, China, October, 23–27 2017.

Automatic Chord Recognition (ACR) Systems

(1) Frame-based/ template-based approach

Realtime Chord Recognition of Musical Sound: a System Using Common Lisp Music

Takuya FUJISHIMA *
CCRMA, Stanford University
Stanford CA 94305, USA
fujishim@ccrma.stanford.edu

Abstract

This paper describes a realtime software system which recognizes musical chords from input sound signals. I designed an algorithm and implemented it in Common Lisp Music/Realtime environment. The system runs on Silicon Graphics workstations and on the Intel PC platform. Experiments verified that the system can recognize chords even in a complex musical environment.

Keywords: music, sound, chord, machine, recognition, realtime, numerical, AI

1 Introduction

Musical chord recognition is a process of identifying specific pitch intervals that combine to produce a functional harmonic palette of western music¹. Traditionally, it is approached as a polyphonic transcription task to identify individual notes and their harmonic context [2] [3], and then follow by a basic inference stage to determine chords (Fig. 1) [4]. This approach suffers from recognition errors at the first stage. The errors result from noises, and from overlapping harmonics of notes in the spectrum. In addition, audio signal processing techniques makes the chord recognition task difficult.



Figure 1: Traditional Approach

Recently, Leman proposed an alternative approach [5]. His "Simple Auditory Model" (SAM for short) avoids the unreliable note identification process by adopting numerical representation and manipulation of intensity maps throughout the analysis. SAM is a robust recognizer of chords. The key of SAM is the use of an intensity map of twelve semitone pitch classes. Derived from a spectrum through straightforward numerical processing, the map preserves more information and achieves more

¹Currently at YAMAHA Corp., Electronic Musical Instrument Division, Japan. E-mail: fujishim@ccrma.stanford.edu

²Human listeners can recognize chords even in a complex musical environment. This paper is beyond the scope of this research.

2 Algorithm

The overview of my chord recognition algorithm is shown in Fig. 2.

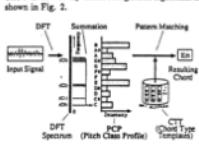


Figure 2: Chord Recognition Algorithm Overview

This algorithm first transforms an input sound to a Discrete Fourier Transform (DFT) spectrum, from which a Pitch Class Profile (PCP) is derived. Then it does pattern matching on the PCP to determine the chord type and root. Two major improvements are introduced to improve the overall performance. Details follow.

2.1 DFT Spectrum

First, this algorithm transforms a fragment of the input sound stream to a DFT spectrum. Let f_s be the sampling frequency and $n(n)$ be a fragment

of the input sound signal with N samples, where $n = 0, 1, \dots, N - 1$ is the sample index. Then the DFT spectrum $X(k)$ is given for $k = 0, 1, \dots, N - 1$ as:

$$X(k) = \sum_{n=0}^{N-1} e^{-j2\pi k n/N} x(n)$$

Note that $X(0), \dots, X(N/2 - 1)$, or the first half, gives the entire spectrum, as $x(n)$ is real.

2.2 Pitch Class Profile

From $X(k)$, this algorithm generates a Pitch Class Profile (PCP). A PCP is similar to the intensity map in Leman's SAM. It is a two-dimensional vector whose elements are the probabilities of the note and semitone pitch classes. For instance, the first PCP element shows how strong C is in total.

Let $PCP(i)$ be a PCP. It is defined for the index $p = 0, 1, \dots, 11$ as

$$PCP(p) = \sum_{k=0}^{N-1} |X(k)|^2$$

where $M(k)$ is a table which maps a spectrum bin index to the PCP index. $M(k)$ can initially be defined as

$$M(k) = \begin{cases} -1 & \text{round}(2\log_2((f_r - \frac{1}{2})/f_s)) \neq 0 \\ 1 & \text{mod}(k, 12) \\ 0 & \text{for } k = 1, 2, \dots, N/12 - 1 \end{cases}$$

where f_s is the reference frequency that falls into $PCP(0)$. The term $(f_r - \frac{1}{2})/f_s$ represents the frequency of the spectrum bin $X(k)$. Note that $PCP(p)$ and $M(k)$ are further modified here for later use (see 2.4).

2.3 Pattern Matching

To improve the overall performance of this algorithm, I introduce the following two heuristics:

1) PCP Smoothing over time: Assuming that a chord usually lasts for several PCP frames, this algorithm derives a smoothed PCP² by averaging successive PCPs to reduce the noise.

2) Chord Change Sensing: A chord change is usually sudden, when the PCP vector changes its direction in the twelve dimension space. Therefore, we can sense chord changes by monitoring the direction of the PCP vector.

Other minor heuristics include:

- PCP preprocessing, including non-linear scaling
- Elimination of less important regions[6] in $M(k)$
- Use of an DFT window
- Silence detection to avoid unnecessary analysis
- Attack detection to avoid noisy PCPs

for which details are not mentioned further here.

2.4 Nearest Neighbor Method

The score to be minimized for a chord type c is given as:

$$\text{Score}_{\text{nearest}} = \sum_{p=0}^{11} (T(p) - PCP(p))^2$$

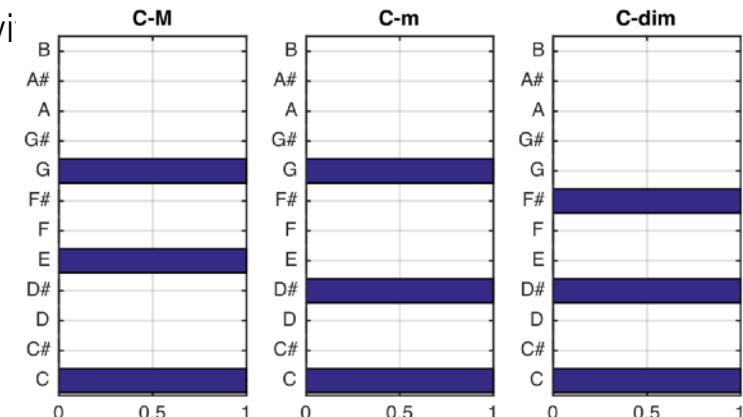
Table 1: Chord Type Chosen

group	chord types
S1	G Major, Gsus, Gmin, G7, Gm7, Gm7sus, G7sus, G7sus4, G7sus2, G7sus9, G7sus11, G7sus13, G7sus15, G7sus17, G7sus19, G7sus21, G7sus23, G7sus25, G7sus27, G7sus29, G7sus31, G7sus33, G7sus35, G7sus37, G7sus39, G7sus41, G7sus43, G7sus45, G7sus47, G7sus49, G7sus51, G7sus53, G7sus55, G7sus57, G7sus59, G7sus61, G7sus63, G7sus65, G7sus67, G7sus69, G7sus71, G7sus73, G7sus75, G7sus77, G7sus79, G7sus81, G7sus83, G7sus85, G7sus87, G7sus89, G7sus91, G7sus93, G7sus95, G7sus97, G7sus99, G7sus101, G7sus103, G7sus105, G7sus107, G7sus109, G7sus111, G7sus113, G7sus115, G7sus117, G7sus119, G7sus121, G7sus123, G7sus125, G7sus127, G7sus129, G7sus131, G7sus133, G7sus135, G7sus137, G7sus139, G7sus141, G7sus143, G7sus145, G7sus147, G7sus149, G7sus151, G7sus153, G7sus155, G7sus157, G7sus159, G7sus161, G7sus163, G7sus165, G7sus167, G7sus169, G7sus171, G7sus173, G7sus175, G7sus177, G7sus179, G7sus181, G7sus183, G7sus185, G7sus187, G7sus189, G7sus191, G7sus193, G7sus195, G7sus197, G7sus199, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7sus299, G7sus201, G7sus203, G7sus205, G7sus207, G7sus209, G7sus211, G7sus213, G7sus215, G7sus217, G7sus219, G7sus221, G7sus223, G7sus225, G7sus227, G7sus229, G7sus231, G7sus233, G7sus235, G7sus237, G7sus239, G7sus241, G7sus243, G7sus245, G7sus247, G7sus249, G7sus251, G7sus253, G7sus255, G7sus257, G7sus259, G7sus261, G7sus263, G7sus265, G7sus267, G7sus269, G7sus271, G7sus273, G7sus275, G7sus277, G7sus279, G7sus281, G7sus283, G7sus285, G7sus287, G7sus289, G7sus291, G7sus293, G7sus295, G7sus297, G7

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) wi
 - 1 if chroma exists in the chord
 - 0 if chroma does not exist in the chord
 - Template $G_a(c)$
 - chord-name, $a \in \{\text{C-M}, \text{C-m}, \text{C\#-M}, \text{C\#-m}, \dots\}$
 - 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



Automatic Chord Recognition (ACR) Systems

(1) Frame-based/ template-based approach

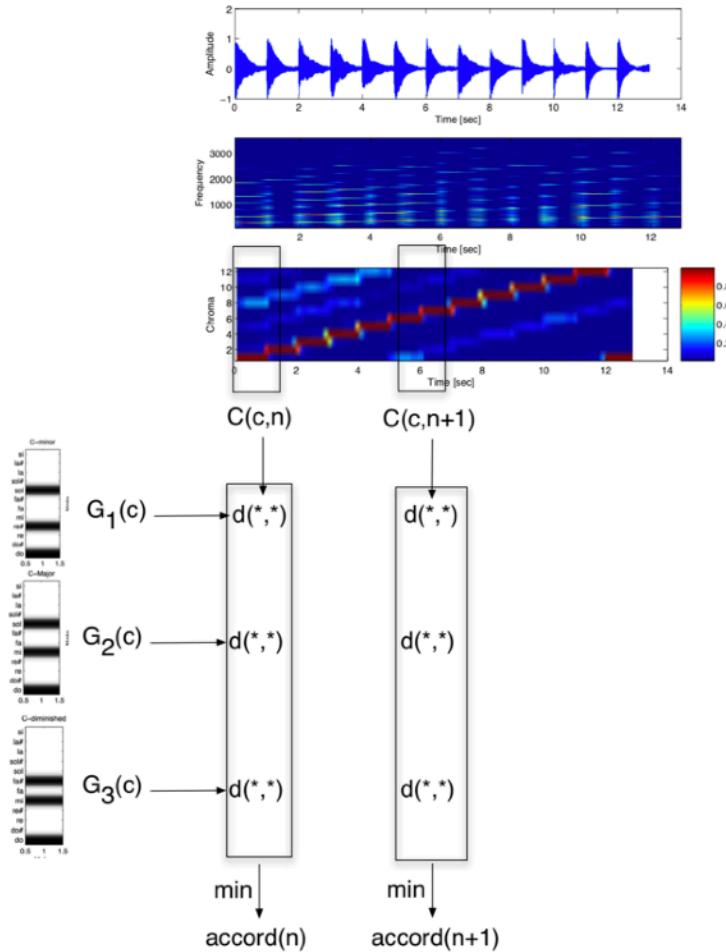
– Frame-based estimation

– 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

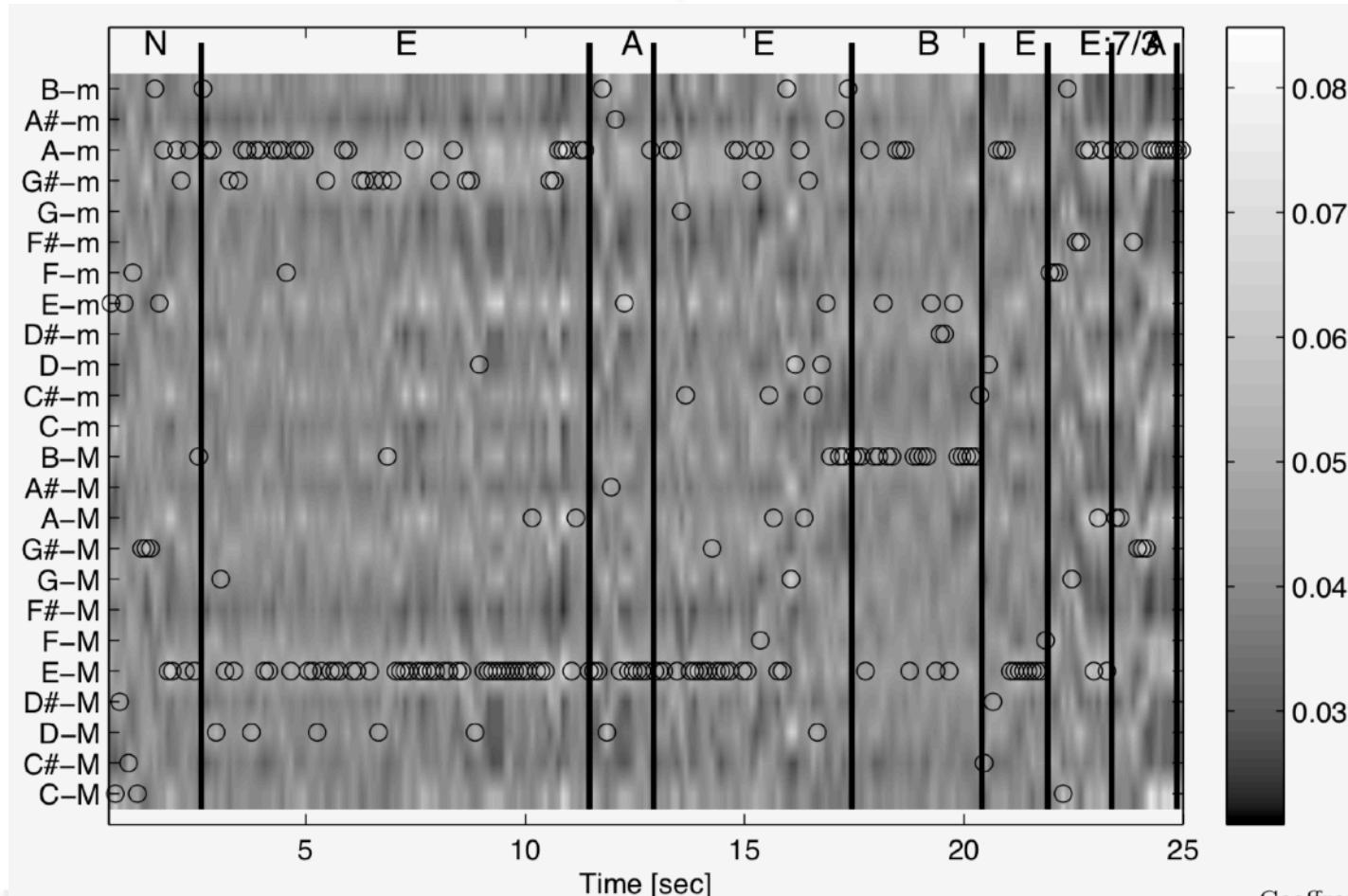
- $d(C(c, n), G_a(c))$ with d
 - Euclidean distance
 - cosine distance
- Choose the chord with the smallest distance
 - $\arg \min_a d(C(c, n), G_a(c))$



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

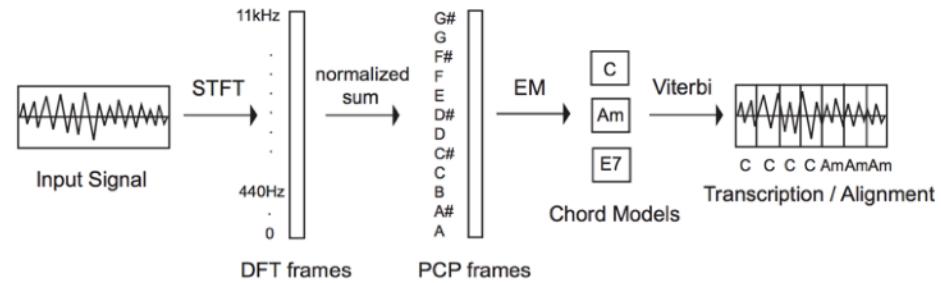
Automatic extraction of content description from commercial audio recordings has become of important applications, from indexing and retrieval of music to novel musicological analyses based on very large corpora of recorded performances. Chord sequences are a description that captures much of the character 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 of a single phonetic state. 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 sequence, our own chord changes — 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 (International Society for Music Information Retrieval), pages 183–189, Baltimore, Maryland, USA, 2003

H. Papadopoulos and G. Peeters. Large-scale study of chord estimation algorithms based on chroma representation. In Proc. of IEEE CBMI (International Workshop on Content-Based Multimedia Indexing), Bordeaux, France, 2007

Automatic Chord Recognition Systems

(2) HMM-based approach

– Defining HMM for chord recognition

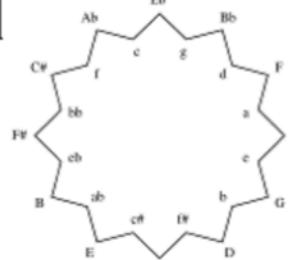
- **Observations O_t**
 - sequence of chroma/PCP vectors over time t
- **States S_1, S_2, \dots, S_N**
 - the chord labels (C-M, C-m, C#-M, ...)
- **Initial state distribution $\pi = \{\pi_i\}$**
 - set to uniform distribution (if we don't have information)
- **Emission probabilities $b_j(O_t) = p(O_t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

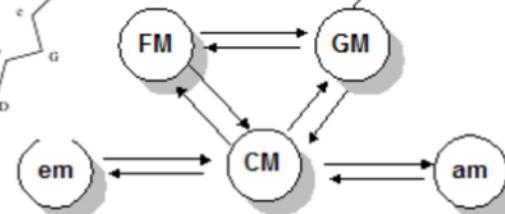
Hidden Markov Model

Modeling
transition
between



...

Modeling
chords



Automatic Chord Recognition Systems

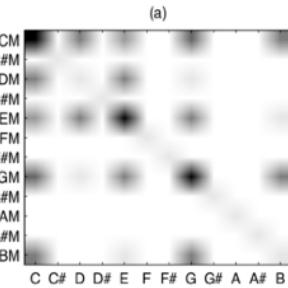
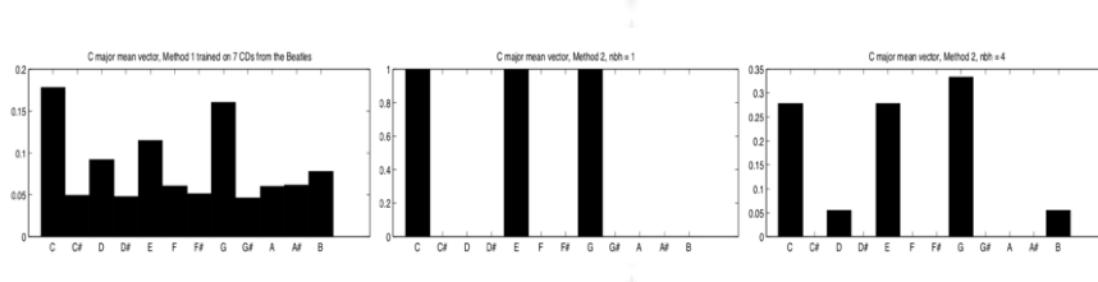
(2) HMM-based approach

– 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 (normalized) distance to chord-templates (see before)
- can be based on manually-tuned statistical models
 - multivariate Gaussian models with parameters μ and Σ
 - 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-u-jah

B₇ | F₇ | C₋₇ F₇ |
B₇ | B₇ | F₇ | A₋₇ D₇ |
G₋₇ | C₇ | F₇ | D₇ | G₋₇ C₇ |
F₇ | B₇ | F₇ | C₋₇ F₇ |
B₇ | B₇ | F₇ | A₋₇ D₇ |
G₋₇ | C₇ | F₇ | D₇ | G₋₇ C₇ |

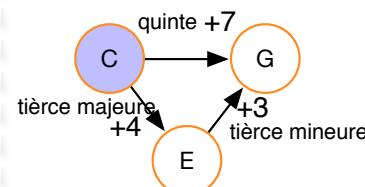
Automatic Chord Recognition Systems

(2) HMM-based approach

- Transition probability between chords:

- can be **based on musical rules**

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

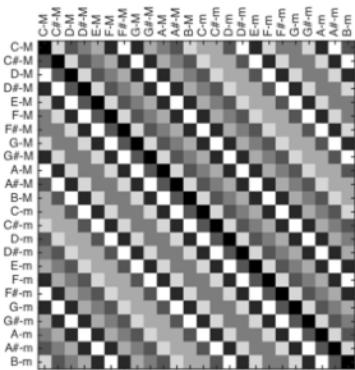
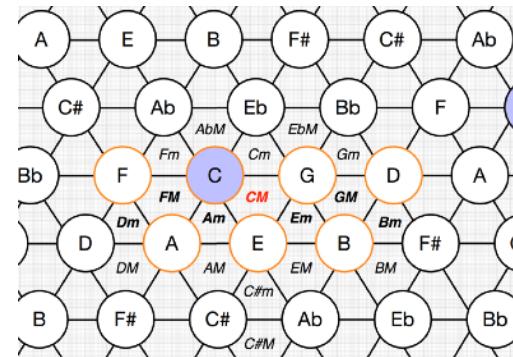
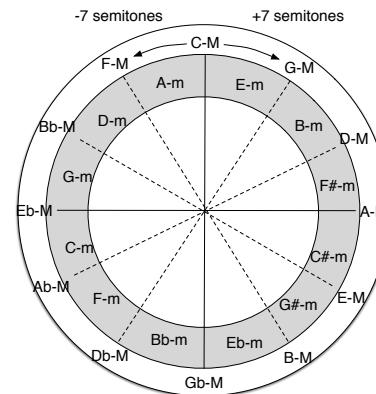


- distance between chords in the Circle-of-Fifths

- séquence de quintes
 - relatifs majeur-mineur
 - Exemple:

G-M vers C-M (consonance),

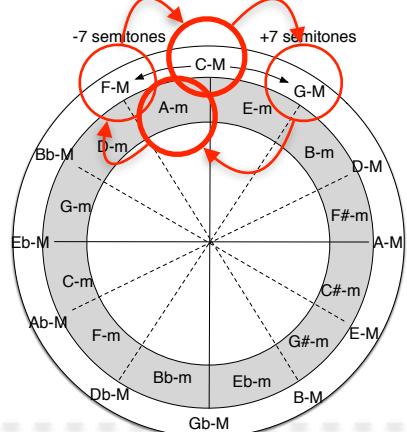
G-M vers Db-M (dissonance)



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 in pop-music:



G7 → Cmaj7

D-7 → G7 → Cmaj7

A-7 → D-7 → G7 → Cmaj7

C7 → F7 → C7 → C7 ||

F7 → F7 → C7 → C7 |

G7 → F7 → C7 → G7 ||

C → G → A- → F

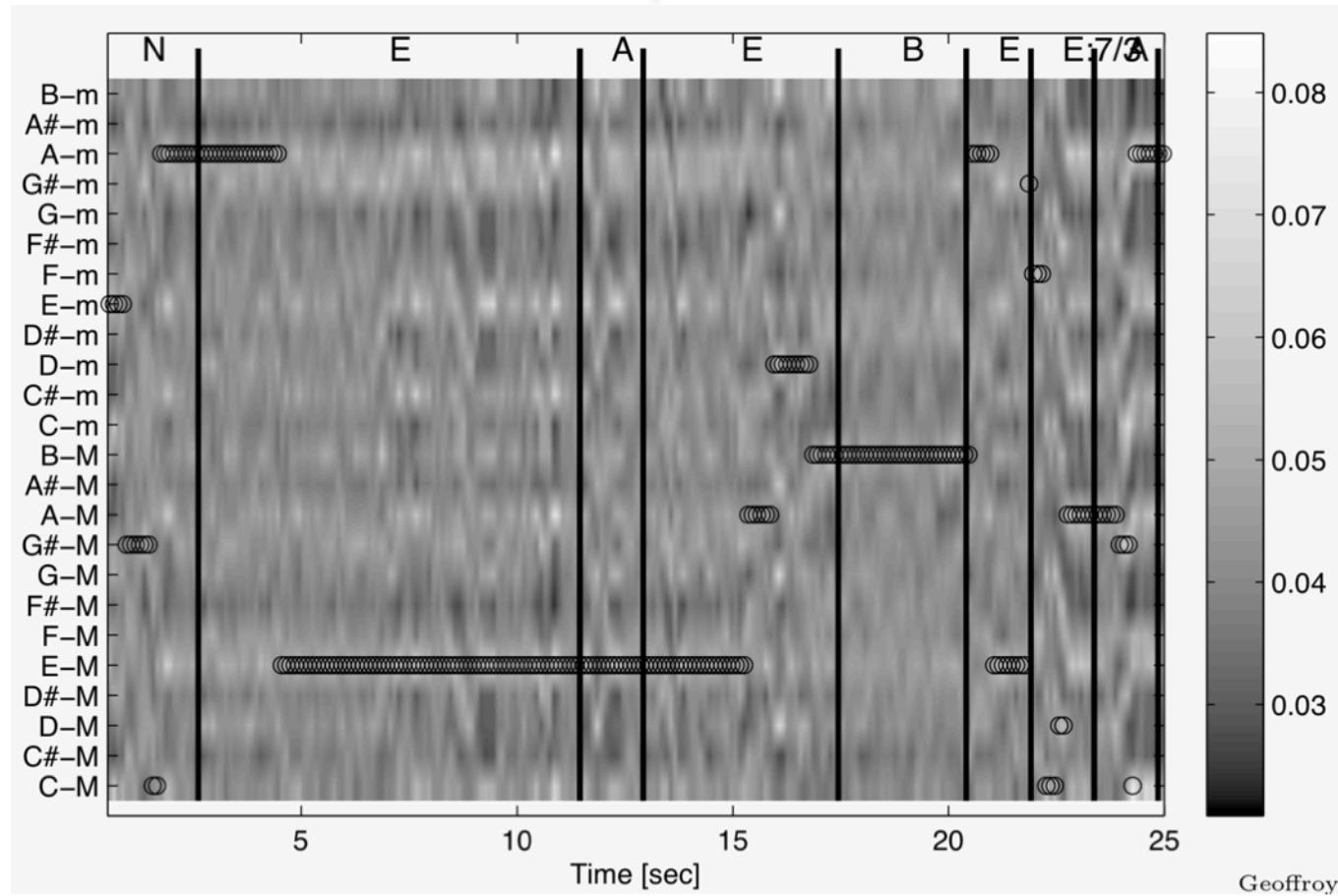
$$A \rightarrow F \rightarrow C \rightarrow G$$



Automatic Chord Recognition Systems

(2) HMM-based approach

– Result example



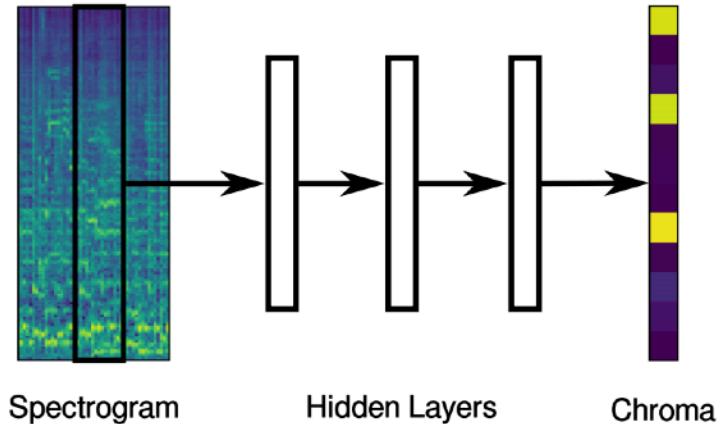
Automatic Chord Recognition (ACR) Systems

(2) Deep-learning-based approach A

- **Goal:**
 - to replace the Chroma/PCP front-end by learned features
 - train a layer-layer MLP to output a ground-truth chroma representation (the one corresponding to the notes of the chord)

– Deep Chroma

- **Evaluation**
 - plug the output to a simple logistic regression to estimate the chord (no post-processing, smoothing)



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

C : standard chroma from CQT

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

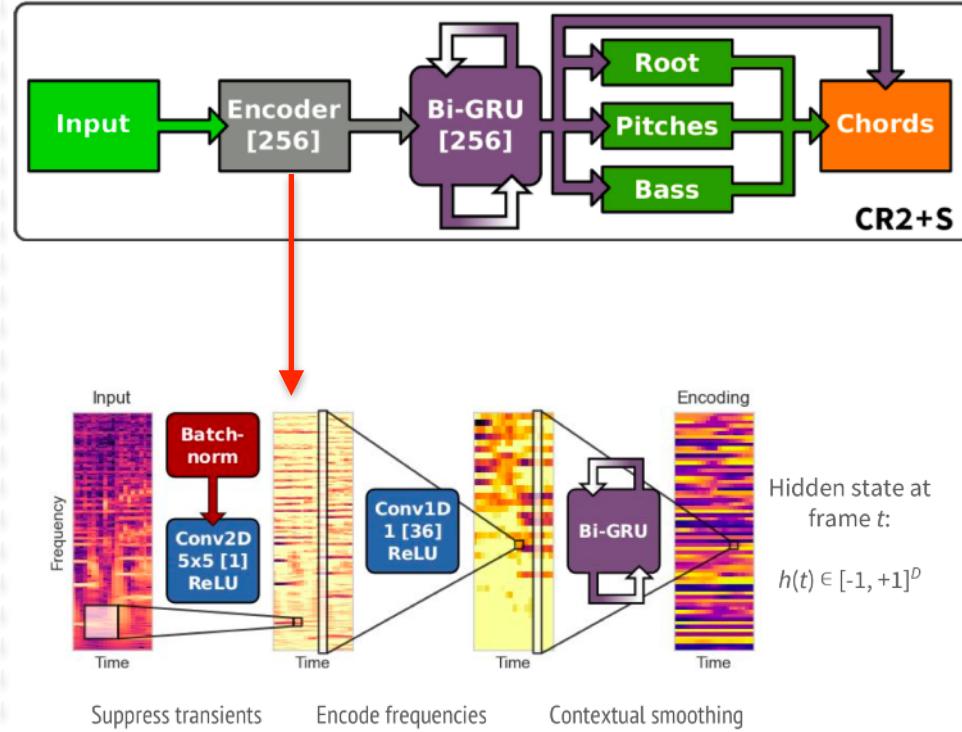
S_{log} : quarter-tone spectrogram

C_D : deep-chroma

Automatic Chord Recognition (ACR) Systems

(2) Deep-learning-based approach B

- Goal 1:
 - End-to-end system



B. McFee and J. P. Bello. Structured training for large-vocabulary chord recognition. In Proc. of ISMIR (International Society for Music Information Retrieval), Suzhou, China, October, 23–27 2017.

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
 - $C:7 = C:\text{maj} + m7$
 - $C:\text{sus4}$ vs. $F:\text{sus2}$
 - Class distribution is non-uniform
 - Rare classes are hard to model
 - Take into account the fact that **some mistakes are better than others**

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



Automatic Chord Recognition (ACR) Systems

(2) Deep-learning-based approach B

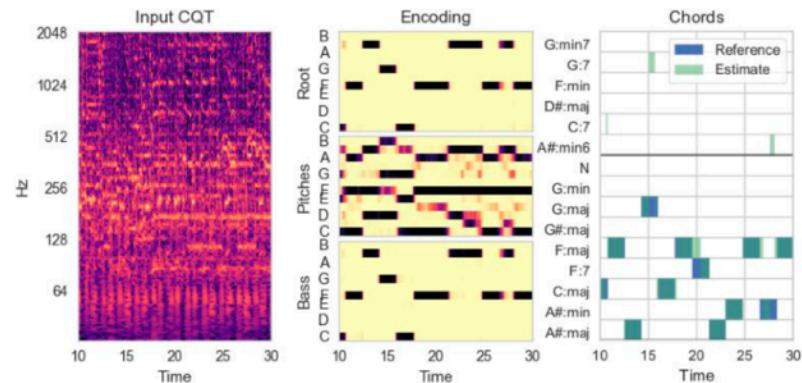
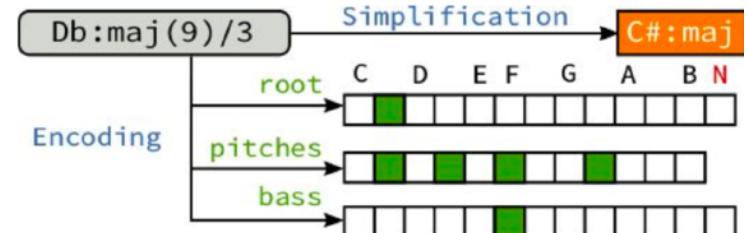
- Idea:
 - Exploit the implicit chord space structure
- Represent chord labels as **binary encodings**
 - Similar chords with different labels will have similar encodings
 - Dissimilar chords will have dissimilar encodings

– Learning problem:

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

– Model architecture

- Input: constant-Q spectral patches
- Per-frame outputs:
 - Root [multiclass, 13]
 - Pitches [multilabel, 12]
 - Bass [multiclass, 13]
 - Chords [multiclass, 170]
- Convolutional-recurrent architecture (encoder-decoder)
- End-to-end training



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 observed

- **8 configurations**

- ± data augmentation
 - ± structured training
 - vs. 2 recurrent layers

- **1217 recordings**

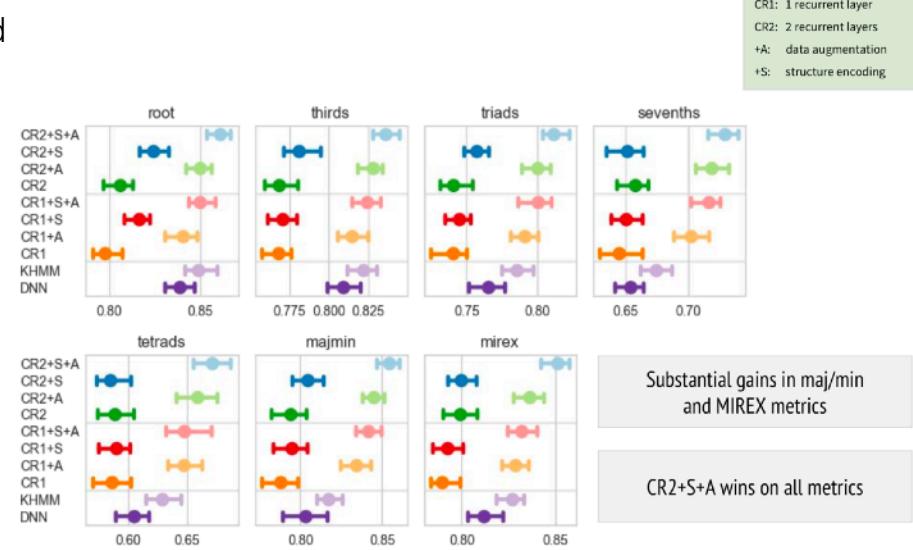
- (Billboard + Isophonics + MARL corpus)
 - 5-fold cross-validation

- **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



Automatic Chord Recognition (ACR)

Evaluation

Automatic Chord Recognition (ACR) Evaluation

Task definition:

- 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	PDF	Chris Cannam  , Matthias Mauch 
JLCX1, JLCX2	PDF	Junyan Jiang  , Ke Chen  , Wei Li  , Guangyu Xia 
SG1	PDF	Franz Strasser  , Stefan Gaser 
FK2	[PDF]	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

(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}

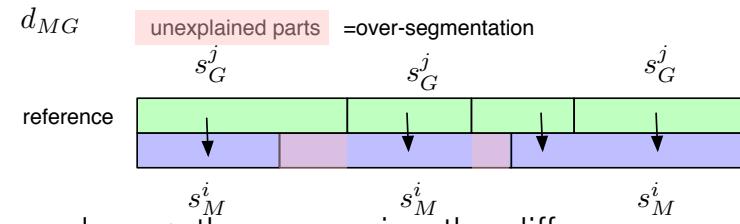
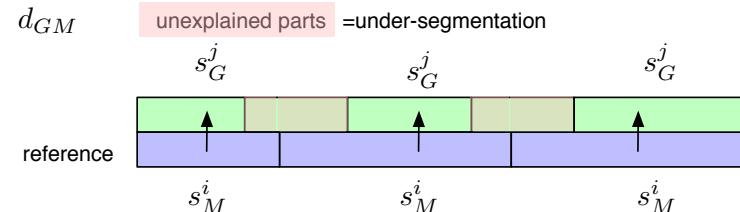
– Vocabulary mapping

- 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}

Automatic Chord Recognition (ACR) Evaluation

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

- 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
 - $d_{GM} = \sum_{S_M^i} \sum_{S_G^k \neq S_G^j} |S_M^i \cap S_G^k|$, 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:** $Q = 1 - \frac{\text{max. of } d_{GM}, d_{MG}}{\text{total duration of song}}$

Chroma/ Pitch Class Profile (PCP)

Chroma/ Pitch Class Profile (PCP)

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

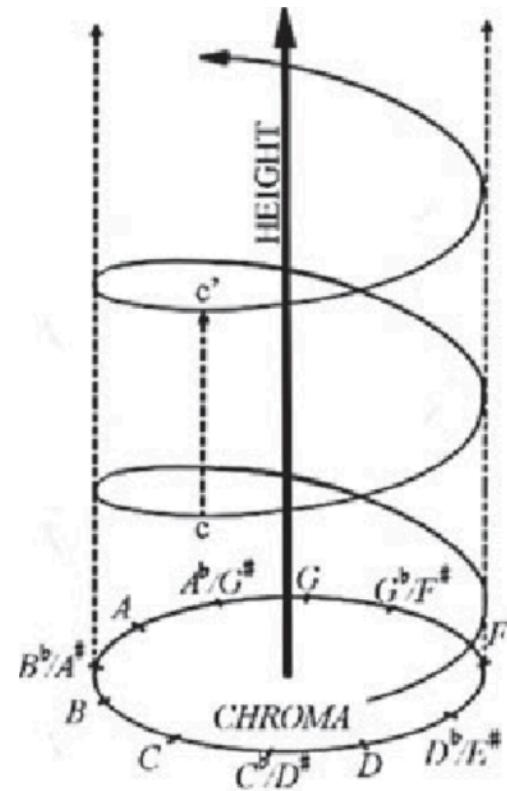
- represents the pitch of a note p as a two-dimensional structure:
- $p = c + o \cdot 12$
 - chroma c (pitch-class)
 - tonal height o (octave number)

- **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

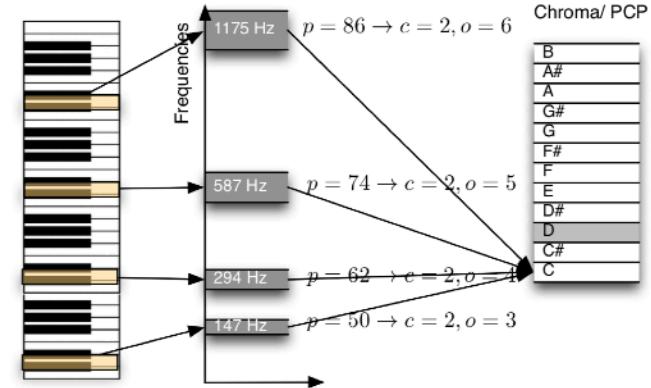
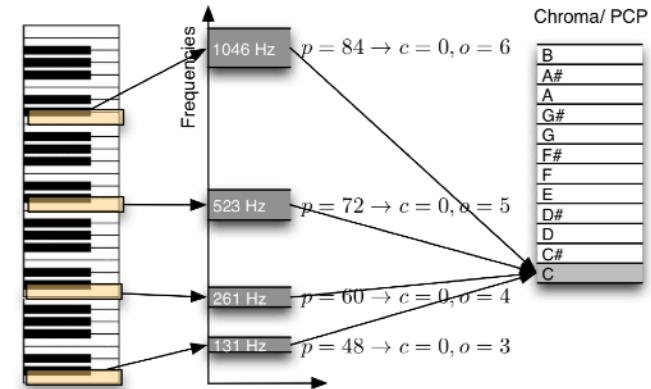


Chroma/ Pitch Class Profile (PCP)

- **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)

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

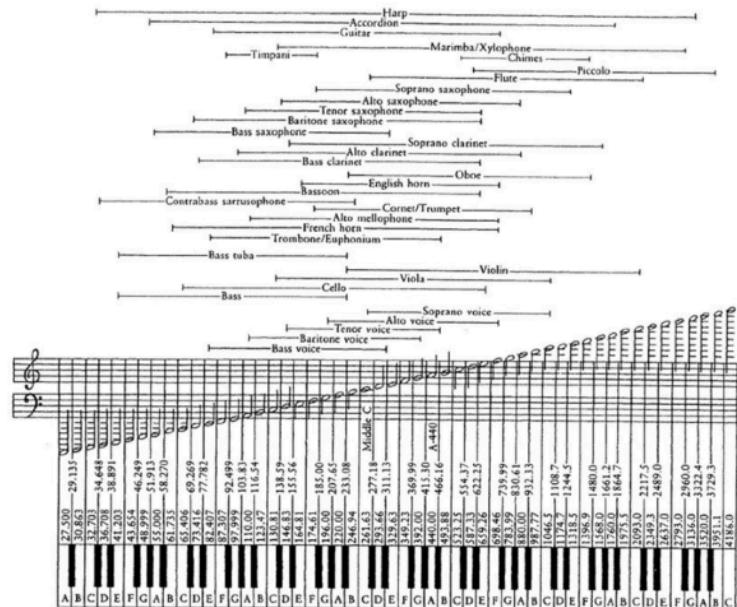
- If f_{min} is 100 Hz

- 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)

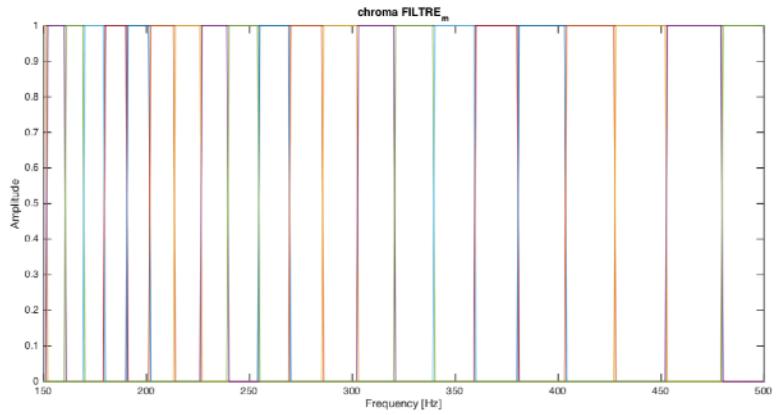
- **Chroma computation** $C(c, n)$

- 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 ?

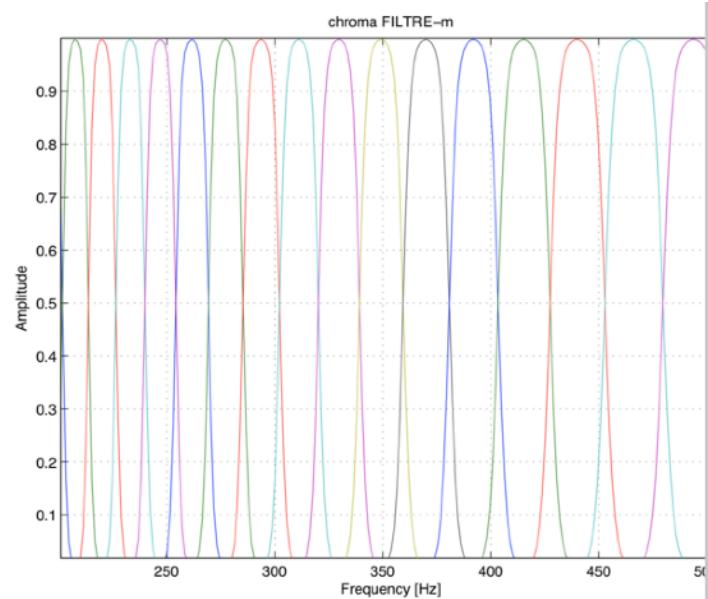
- A frequency f_k of the DFT only contributes to the closest pitch
- Example:
 - the energy at $f_k = 452 \text{ Hz}$ ($p(f_k) = 69.4658$) only contributes to the pitch $p=69$ ($c=10$)
 - while $f_k = 453 \text{ Hz}$ ($p(f_k) = 69.5041$) to $p=70$ ($c=11$).
- Creation of bank of filters $H_{p'}$ centered 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 $H_{n'}$:

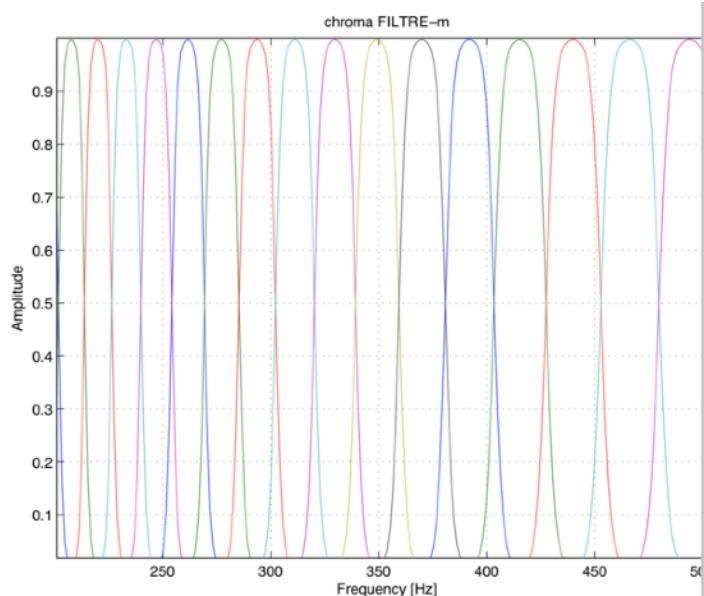
$$P(p') = \sum_{f_k} H_{p'}(f_k)A(f_k)$$

- 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)**

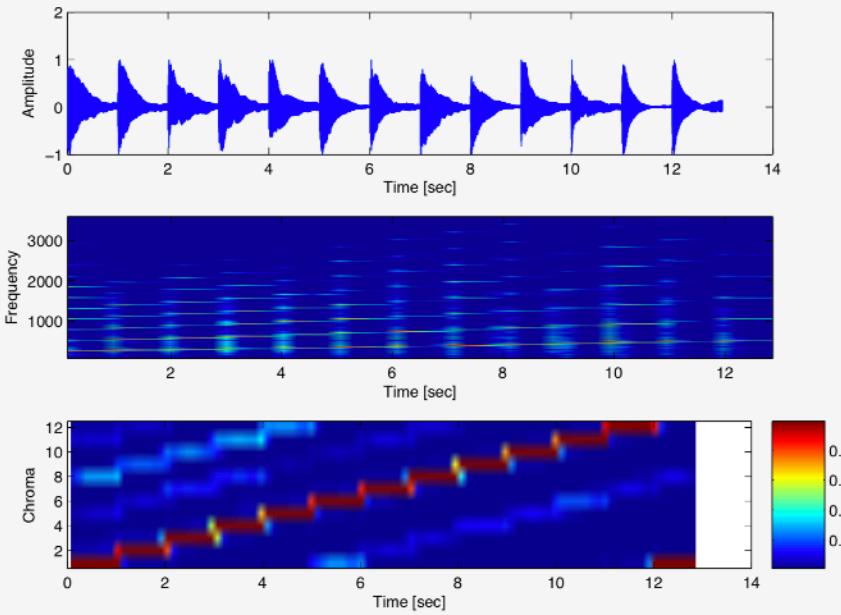
- Presence of the upper harmonics of each note
 - In practice, for a given note C we don't have $[1,0,0,0,0,0,0,0,0,0]$
 - but rather $[a_1 + a_2 + a_4, 0, 0, 0, a_5, 0, 0, a_4, 0, 0, 0, 0]$
 - Influence of the spectral envelope

Pitch	Harmonic	Frequency f_μ	MIDI-scale m_μ	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 ($\simeq e$)

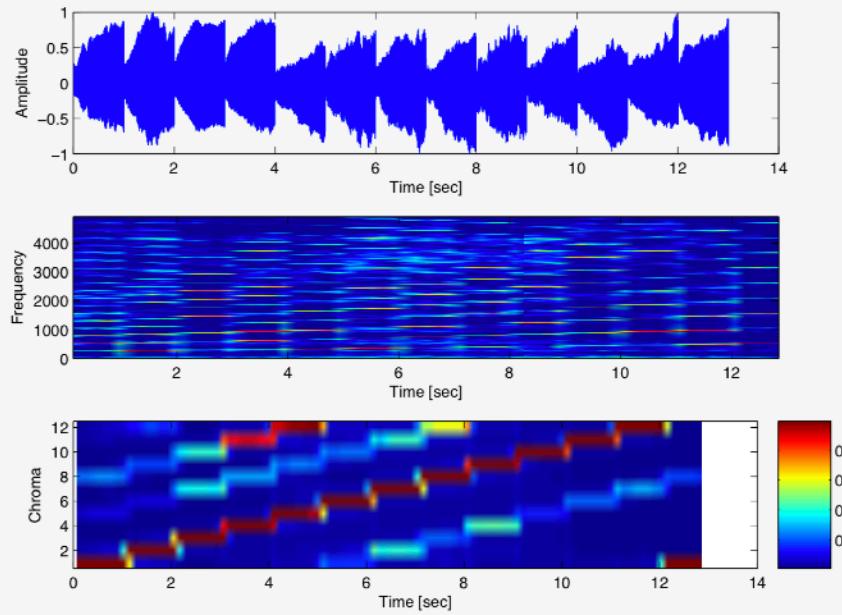
Chroma/ Pitch Class Profile (PCP)

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

Exemple piano

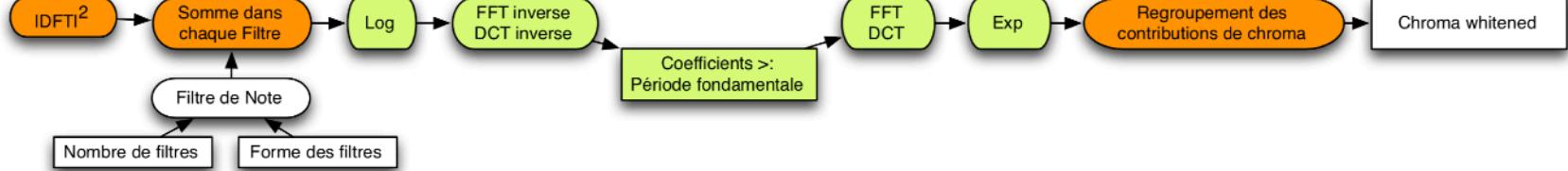
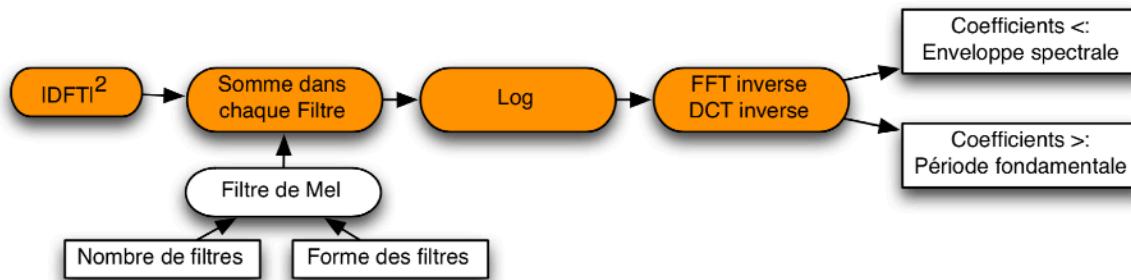
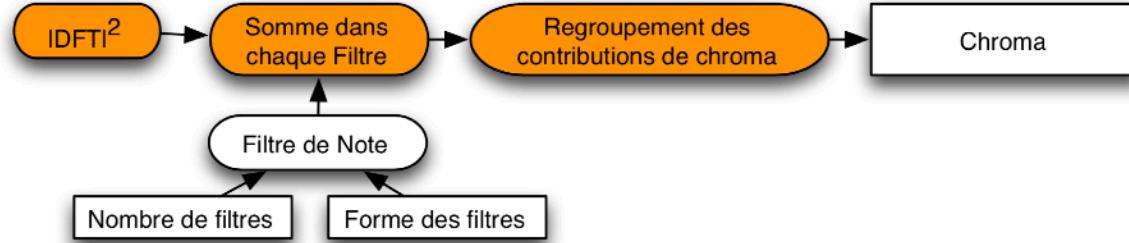


Exemple violon



Chroma/ Pitch Class Profile (PCP)

- Variation of the computation: whitening



Transformée à Q constant

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** : $B_w = \frac{C_w}{L}$

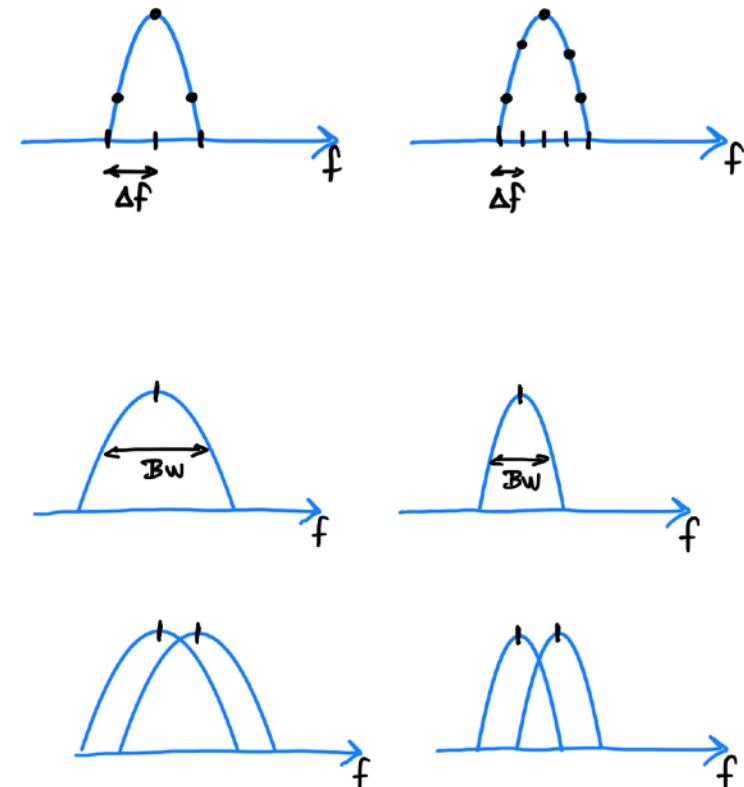
- it describes the ability to discriminate (separate in the spectrum) two adjacent simultaneous frequencies

- Warning :

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

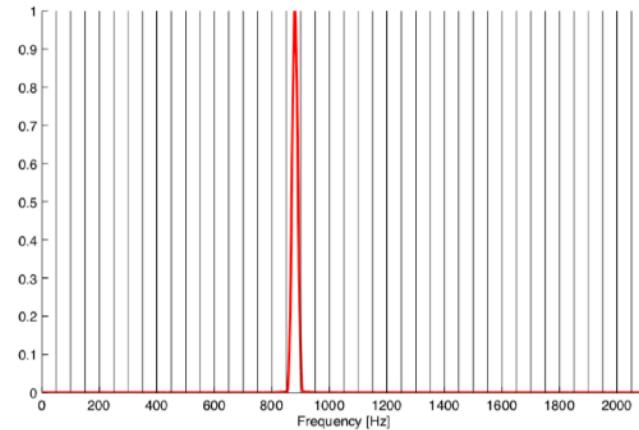
- In a DFT:

- Spectral precision and resolution are constant over frequencies

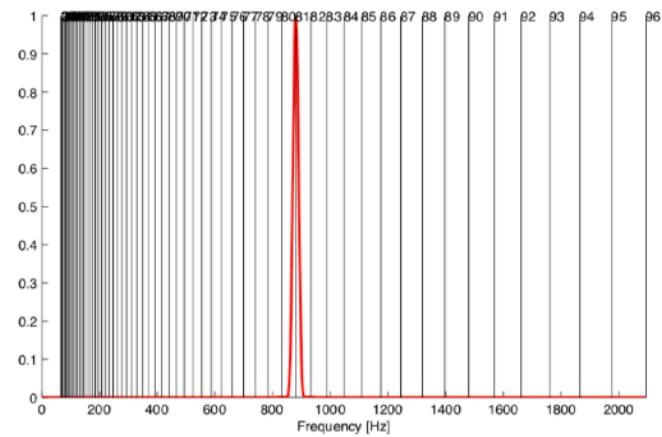


Constant-Q-Transform (CQT)

- 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

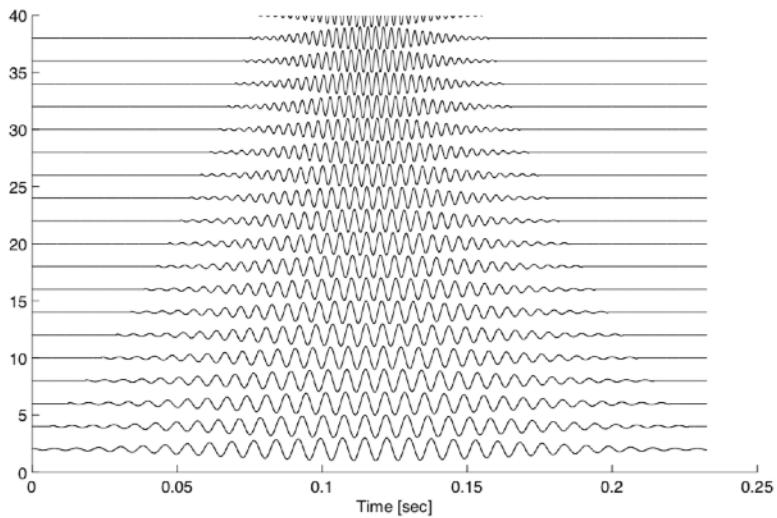


Espacement logarithmique des hauteurs de notes

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
Physics Department, Rensselaer Polytechnic Institute, Troy, New York 12180

Michael S. Puckette
Computer Music Technology Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139

(Received 2 February 1992; revised 10 August 1992; accepted 18 August 1992)

Abstract: An algorithm for calculating a discrete Fourier transform (DFT) into a constant Q transform, where Q is the ratio of center frequency to bandwidth, has been derived. This transform is a generalization of the constant M -band transform, which is a sum of several overlapping windows of length M . The new transform is a sum of several bands of length L , where L is proportional to Q . The computation of the transform is based on a few multiples of the number of samples that are non-overlapping. The computation of the transform requires a small amount of memory. In effect, this method makes it much easier to calculate a constant Q transform than to calculate a constant M -band transform or a short-time Fourier transform (STFT). Graphical examples of the application of the calculation to musical signals are given for sounds produced by a clarinet and a violin.

PACS numbers: 43.60.Dv, 43.75.Fy, 43.75.Dm, 43.75.Ef

I. THEORY

In some cases, such as that of musical signals, a constant Q transform gives a better representation of spectral data than a more commonly used Fourier transform. A recent extension of the constant M -band transform to constant Q transforms has been described by G. Libesman and M. Puckette (1991). The most popular of these transforms is a "wavelet transform," which is a sum of several overlapping windows of length M .

The constant Q transform is a sum of several overlapping windows of length L , where L is proportional to Q . This transform is also called a "constant Q Fourier transform."

We have calculated a constant Q transform based on a DFT, which is a sum of several overlapping windows of length N .

The DFT is calculated using a standard FFT program, and the entire calculation takes only slightly longer than the calculation of a constant M -band transform. The computation of the transform is based on a direct convolution method applied in the time domain.

We have calculated a constant Q transform based on a DFT, which is a sum of several overlapping windows of length N .

The DFT is calculated using a standard FFT program, and the entire calculation takes only slightly longer than the calculation of a constant M -band transform. The computation of the transform is based on a direct convolution method applied in the time domain.

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

$$X^M[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}, \quad (1)$$

where $X^M[k_m]$ is the k_m component of the constant M -band transform, we find that the m th component of each value of $x_m[n]e^{-j2\pi n k_m/N}$ is a window function of length M centered at $n=k_m$. The expression gives the effect of the filter for center frequency f_m .

In constant Q transforms the center frequencies are proportional to the frequency of the signal. The transform is often based on the frequencies of the equal tempered scale with

$$x_m = (2\pi f_m)^{-1/2} \sin(2\pi f_m t_m), \quad (2)$$

for sensible spacing.

f_m is defined as $f_m = f_0 / Q^m$, where f_0 denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (3)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (4)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (5)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (6)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (7)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (8)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (9)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (10)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (11)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (12)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (13)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (14)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (15)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (16)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (17)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (18)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (19)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (20)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (21)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (22)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (23)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (24)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (25)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (26)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (27)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (28)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (29)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (30)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (31)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (32)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (33)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (34)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (35)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (36)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (37)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (38)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (39)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (40)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (41)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (42)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (43)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (44)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (45)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (46)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (47)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (48)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (49)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (50)$$

For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

$$K[k_m] = \sum_{n=0}^{N-1} x[n]e^{-j2\pi n k_m/N}. \quad (51)$$

We will note that $K[k_m]$ is the frequency domain side of the transform, and $X^M[k_m]$ is the frequency domain side of the transform.

We have used a Hamming window:

$$x[n] = 0.54 - 0.46 \cos(2\pi n / N). \quad (52)$$

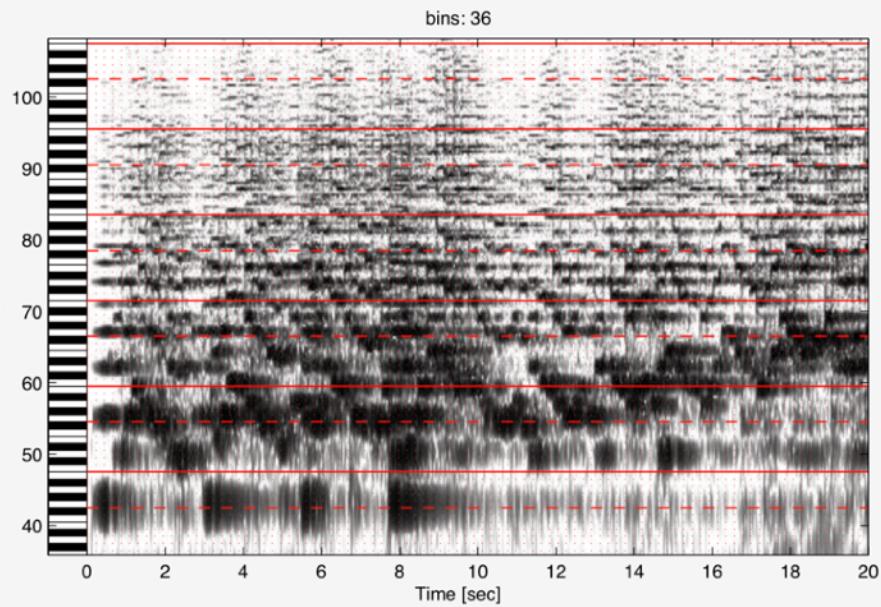
For sensible spacing.

Q is defined as $Q = f_0 / B$, where B denotes bandwidth and f denotes frequency. In the case of the filter defined in Eq. (1), this bandwidth depends on the filter

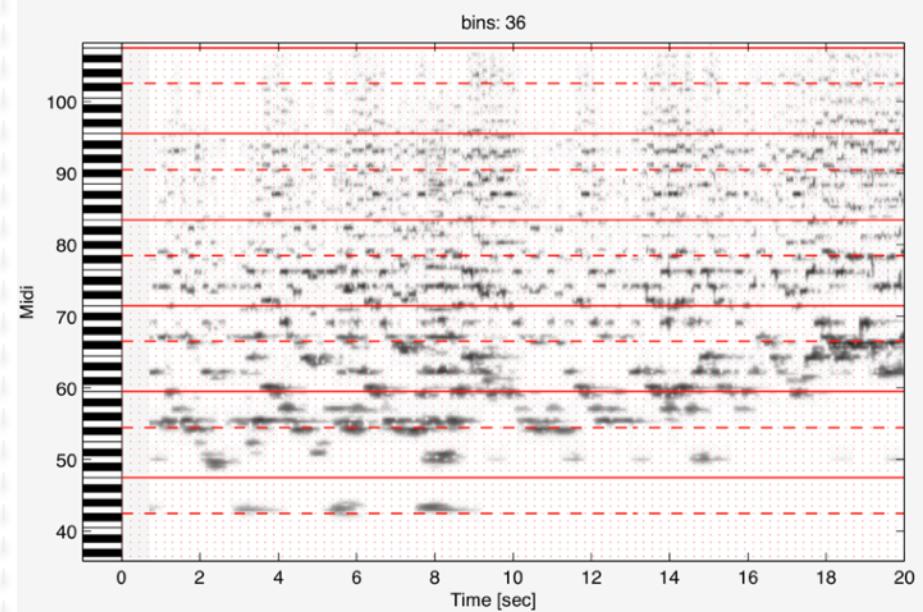
$$K[k_m] = \sum_{n=0}^{N-1}$$

Constant-Q-Transform (CQT)

Example (using DFT)

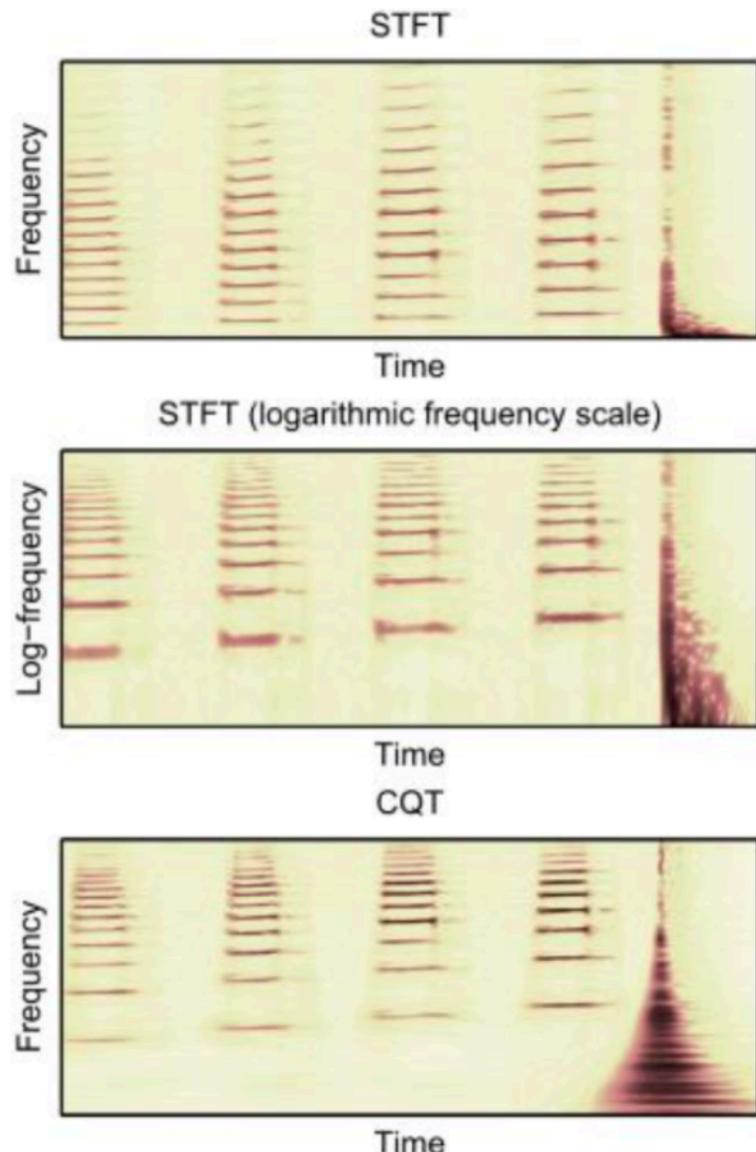


Example (using the CQT)



Constant-Q-Transform (CQT)

- In the Constant-Q-Transform (CQT):
 - A pitch difference corresponds to a translation over the (log) frequency axis



Hidden Markov Model (HMM)

Hidden Markov Model (HMM)

- Modèle de Markov
 - Andreï A. Markov (1856-1922): un mathématicien Russe
- Chaîne de Markov:
 - un processus stochastique à **temps discret** t pouvant être dans des **états discrets** $S_i, i \in \{1, \dots, I\}$
 - q_t la valeur de l'état à l'instant t
- Chaîne de Markov d'ordre 1:
 - la prédiction de l'état actuel ne dépend que de l'instant précédent.
 - $p(q_t | q_{t-1}, q_{t-2} \dots q_0) = p(q_t | q_{t-1})$

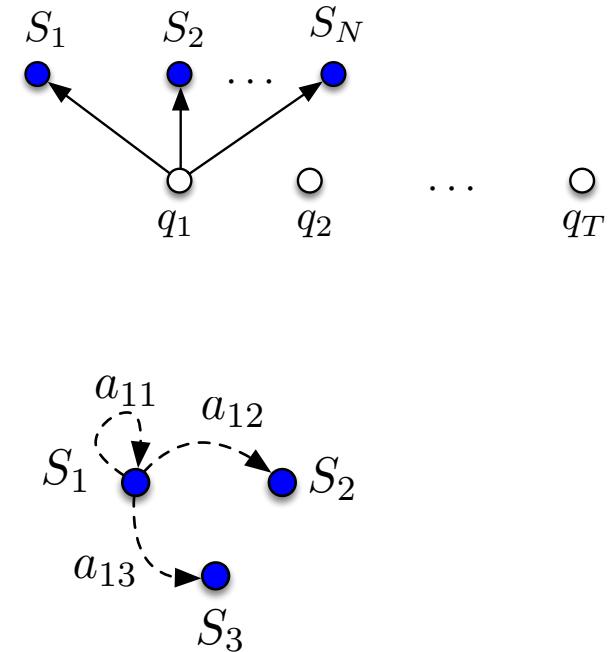


source: [https://fr.wikipedia.org/wiki/Andre%C3%A9_Markov_\(math%C3%A9maticien\)](https://fr.wikipedia.org/wiki/Andre%C3%A9_Markov_(math%C3%A9maticien))

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 $Q = q_1, q_2, \dots, q_T$
- Discrete first order Markov chain
 - $P(q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots) = P(q_t = S_j | q_{t-1} = S_i)$
- **Transition probability** 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
 - $\pi_i = P(q_1 = S_i)$

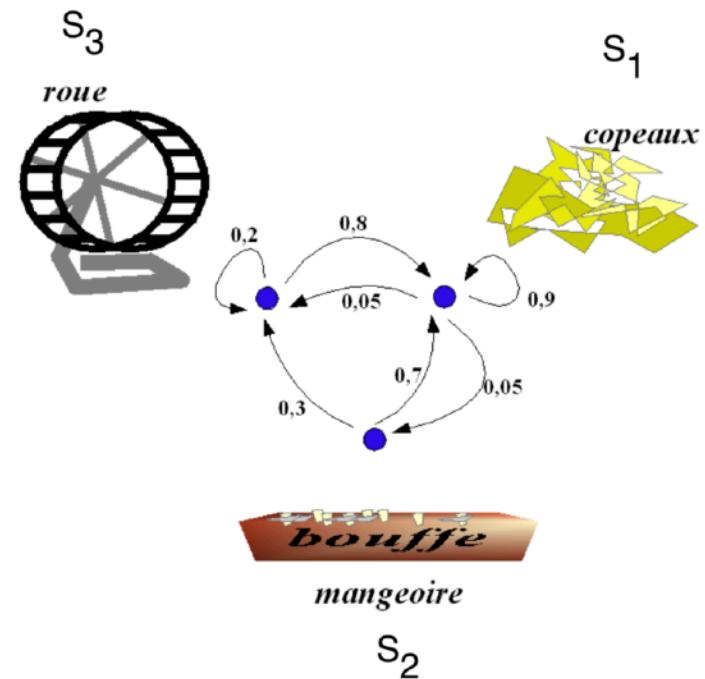


Hidden Markov Model (HMM)

Observed Markov model example

- Doudou le hamster a 3 états dans sa journée:
 - il dort:
 - il est dans l'état S_1 (copeaux)
 - il mange:
 - il est dans l'état S_2 (mangeoire)
 - il fait du sport:
 - il est dans l'état S_3 (roue)
- On peut représenter la succession de ces états par une **matrice de transition** $A = (a_{ij})$ entre états

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



Modèle de Markov d'une journée de Doudou le hamster

source: https://fr.wikipedia.org/wiki/Cha%C3%ABne_de_Markov

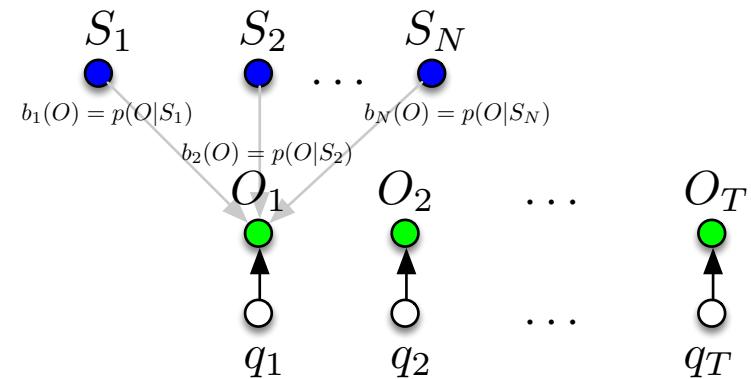
Hidden Markov Model (HMM)

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** $O = O_1, O_2, \dots, O_T$

– 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

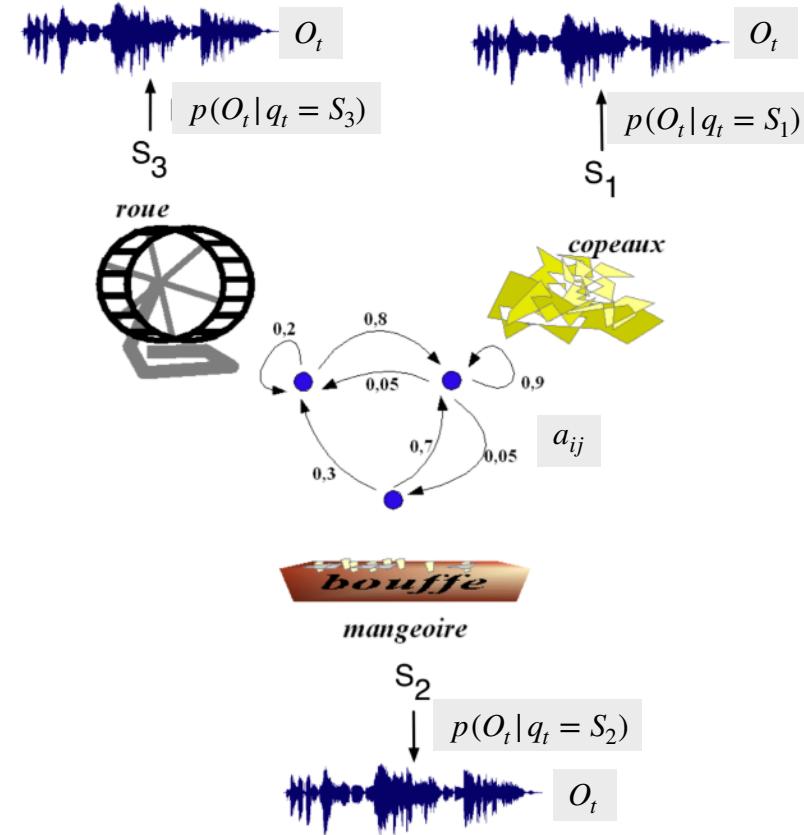
- Gaussian Mixture Model
- $b_j(O) = \sum_m c_{jm} \mathcal{N}(O, \mu_{jm}, \Sigma_{jm})$

- The observation is a probabilistic function of the state S_j
- We denote by $\lambda = \{A, B, \pi\}$ the set of elements defining an HMM

Hidden Markov Model (HMM)

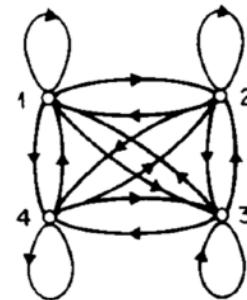
Hidden Markov model example

- Dans un modèle de Markov caché on observe pas directement les états $Q = q_1, q_2, \dots, Q_T$, ils sont "cachés"
- on observe une émission des états q_t
 - exemple: on observe le bruit $O = O_1, O_2, \dots, O_T$ que fait Doudou le hamster
- Pour chaque état, nous pouvons cependant définir la
 - probabilité d'émission de O étant donné l'état S_i :
 $b_j(O_t) = p(O_t | q_t = S_i)$

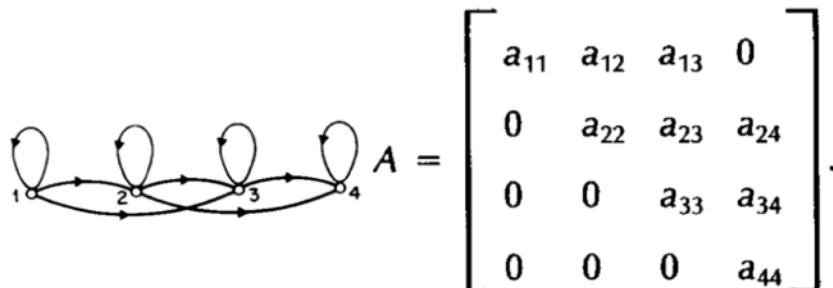


Différentes Topologies des modèles de Markov

- 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)

– Etant donné

- une suite d'observation $O = O_1, O_2, \dots, O_T$
- un modèle $\lambda = \{A, B, \pi\}$

– Quelle est la probabilité que ce modèle ait généré O ?

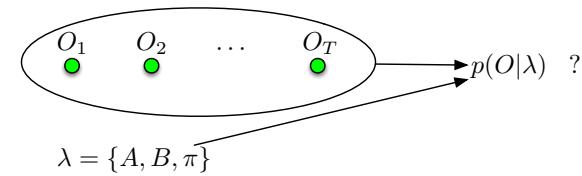
- $P(O|\lambda)$?

– Exemple d'application:

- comment déterminer si une séquence d'observation du son O correspond
 - au modèle λ_1 {dormir/manger/exercice} de Doudou le hamster ou
 - au modèle λ_2 {dormir/manger/travailler} de Bill le salarié

– **Problème:**

- On ne sait pas par quel états $Q = q_1, q_2, \dots, q_T$ le modèle est passé
 - On doit tous les considérer

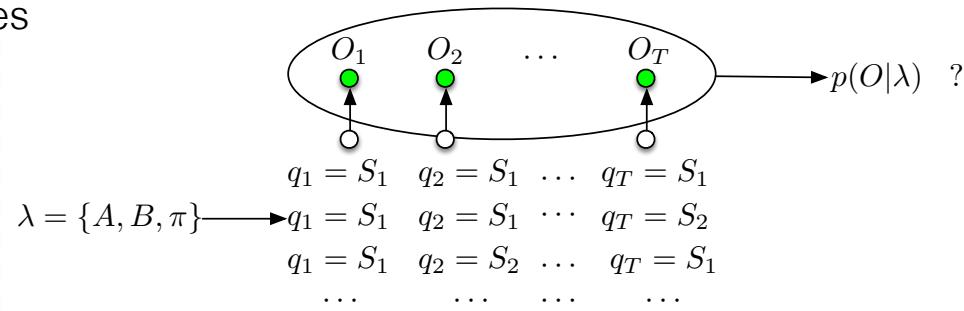


Hidden Markov Model (HMM)

Problem 1. Likelihood of a model (Forward algorithm)

Méthode 1 (exhaustive)

- Nous devons énumérer toutes les séquences d'états de longueur T : $Q = q_1, q_2, \dots, q_T$



- Pour chaque séquence Q , nous estimons

- la probabilité de ces observations

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

- la probabilité de cette séquence est

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

- la probabilité jointe de O et Q

$$P(O, Q|\lambda) = P(O|Q, \lambda)P(Q|\lambda)$$

- Nous devons finalement sommer pour toutes les séquences Q

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

– Cout de calcul énorme !!!

Hidden Markov Model (HMM)

Problem 1. Likelihood of a model (Forward algorithm)

Méthode 2 (optimisée): Forward algorithm

– Definition: **Forward variable**

- $\alpha_t(j) = P(O_1, O_2, \dots, O_t, q_t = S_j | \lambda)$
- 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 λ

• Initialisation

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

• Induction

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

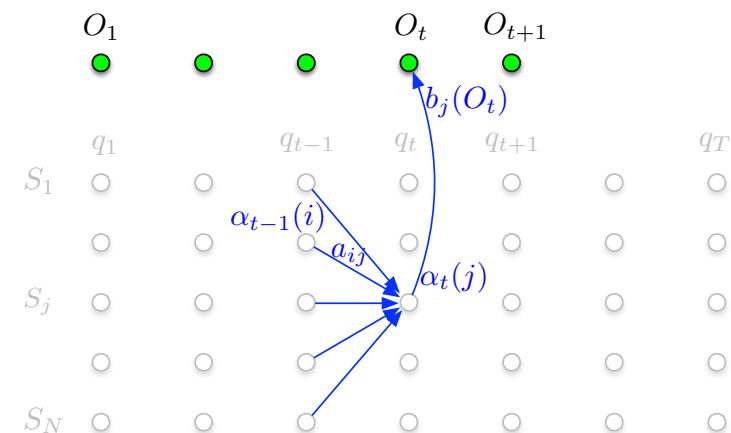
• Terminaison

$$-\sum_{i=1}^N \alpha_T(i)$$

– Preuve:

$$\alpha_T(i) = P(O_1, O_2, \dots, O_T, q_t = S_i | \lambda)$$

$$\sum_i \alpha_T(i) = \sum_i P(O_1, O_2, \dots, O_T, q_t = S_i | \lambda)$$
$$= P(O_1, O_2, \dots, O_T | \lambda)$$



Hidden Markov Model (HMM)

Problem 1. Likelihood of a model (Forward algorithm)

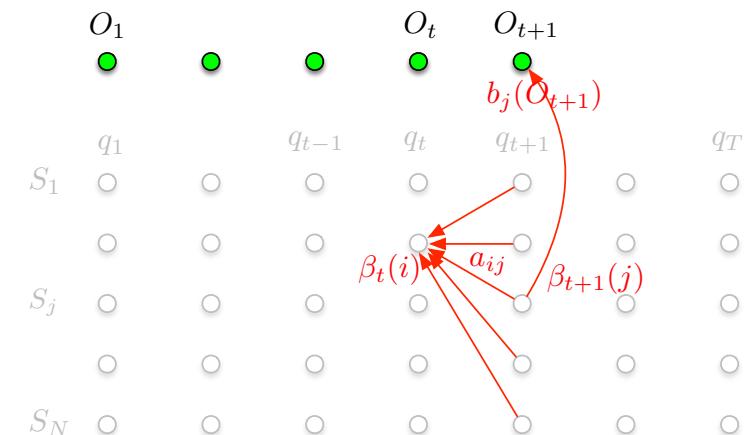
Méthode 2 (optimisée): revisited

- We could also have used the reverse recursion

- Definition: **Backward variable**

- $\beta_t(i) = P(O_{t+1}, O_{t+2}, \dots, O_T | q_t = S_i, \lambda)$
- probability of the partial observation sequence from $t + 1$ to end
 - given state S_i at time t , and the model λ

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

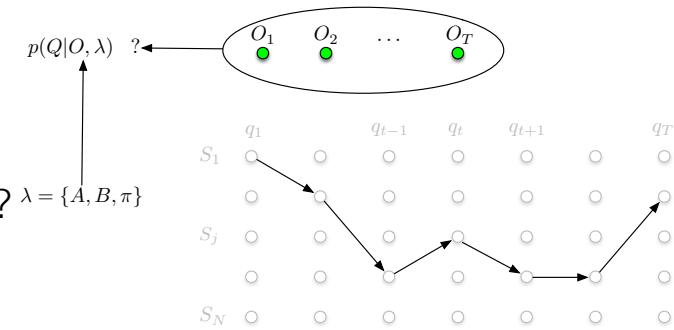


Hidden Markov Model (HMM)

The three basic problems for HMMs

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

- Etant donné
 - une suite d'observation $O = O_1, O_2, \dots, O_T$
 - un modèle $\lambda = \{A, B, \pi\}$
- Quelle est la suite d'état $Q = q_1, q_2, \dots, q_T$ correspondant ?
 - $Q^* = \arg \max_Q P(Q | O, \lambda)$?
- Exemple d'application:
 - si on observe la séquence de son O de Doudou et étant donné son modèle λ {dormir/manger/exercice}, quelle est la séquence d'activités Q de Doudou ?
- **Problème:**
 - Recherche de la suite Q la plus optimale
 - Plusieurs définition d'optimalité possible



Hidden Markov Model (HMM)

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

Optimalité 1) choisir les q_t qui sont individuellement les plus probables

– Définition

- $\gamma_t(i) = P(q_t = S_i | O, \lambda)$
- probability of being in state S_i at time t given the observation sequence O and the model λ

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

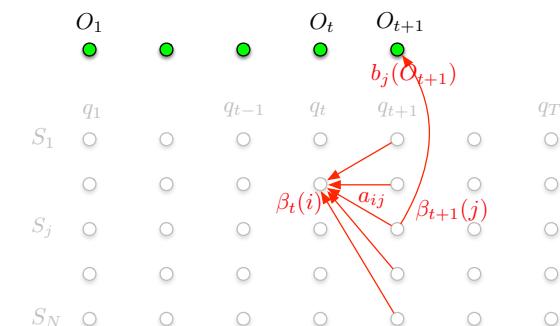
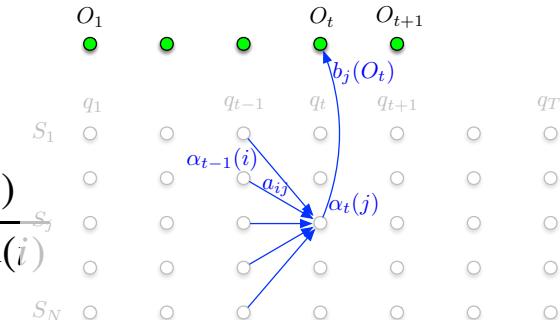
$$\gamma_t(i) = \frac{p(q_t = S_i, O | \lambda)}{p(O | \lambda)} = \frac{\alpha_t(i)\beta_t(i)}{\sum_q p(q, O | \lambda)} = \frac{\alpha_t(i)\beta_t(i)}{\sum_i \alpha_t(i)\beta_t(i)}$$

– Choix de $q_t = \arg \max_i [\gamma_t(i)]$

- problème: ne prend pas en compte les probabilités de transitions (what about if $a_{ij} = 0$?)

• Rappel:

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



Hidden Markov Model (HMM)

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

Optimalité 2) choisir la séquence Q qui est globalement la plus probable

- Programmation dynamique

- Définition: $\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2 \dots q_t = S_i, O_1, O_2 \dots O_t | \lambda)$
- highest probability of the partial alignment starting at 1 and ending at t in state S_i
- On peut calculer $\delta_t(\cdot)$ à partir de $\delta_{t-1}(\cdot)$ comme
 - $\delta_t(j) = \max_i [\delta_{t-1}(i) a_{ij}] b_j(O_t)$

- Algorithme de Viterbi

- Forward pass

- Initialisation:

$$\delta_1(i) = \pi_i b_i(O_1) \quad \psi_1(i) = 0$$

- Récursion

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

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

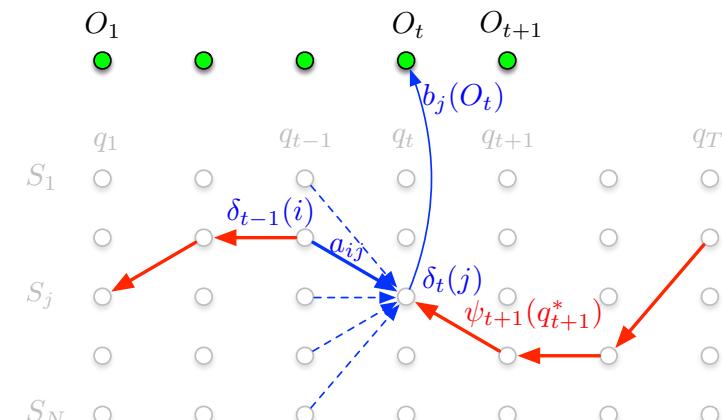
- Terminaison

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

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

- Backward pass

- $q_t^* = \psi_{t+1}(q_{t+1}^*) \quad t = T-1, T-2, \dots, 1$

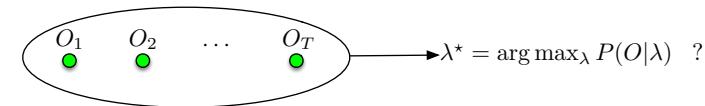


Hidden Markov Model (HMM)

The three basic problems for HMMs

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

- Etant donné
 - une/des suites d'observations $O = O_1, O_2, \dots, O_T$
- Quel est le modèle λ^* qui maximise la vraisemblance des observations:
 - $\lambda^* = \arg \max_{\lambda} P(O | \lambda)$
 - i.e. trouver les paramètres $\lambda = \{A; B, \pi\}$ qui maximise la probabilité de la séquence d'observations O
- Exemple d'application:
 - comment déterminer les paramètres λ du modèle de Doudou le hamster ?



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

Baum-Welch algorithm

- EM expectation-maximization algorithm

- Definition

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

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

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

- On a aussi

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

- Rappel:**

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

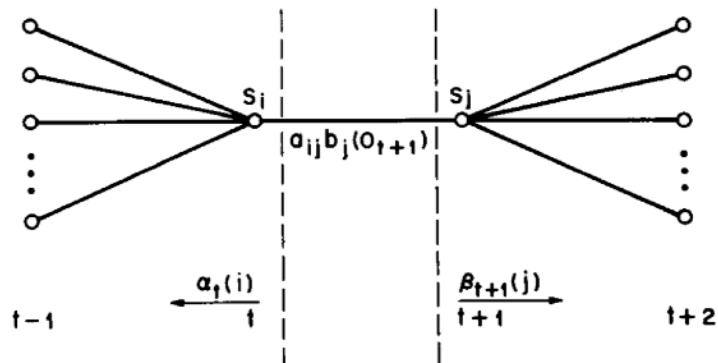


Fig. 6. Illustration of the sequence of operations required for the computation of the joint event that the system is in state S_i at time t and state S_j at time $t + 1$.

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

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

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

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

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

– HMM parameters constraints

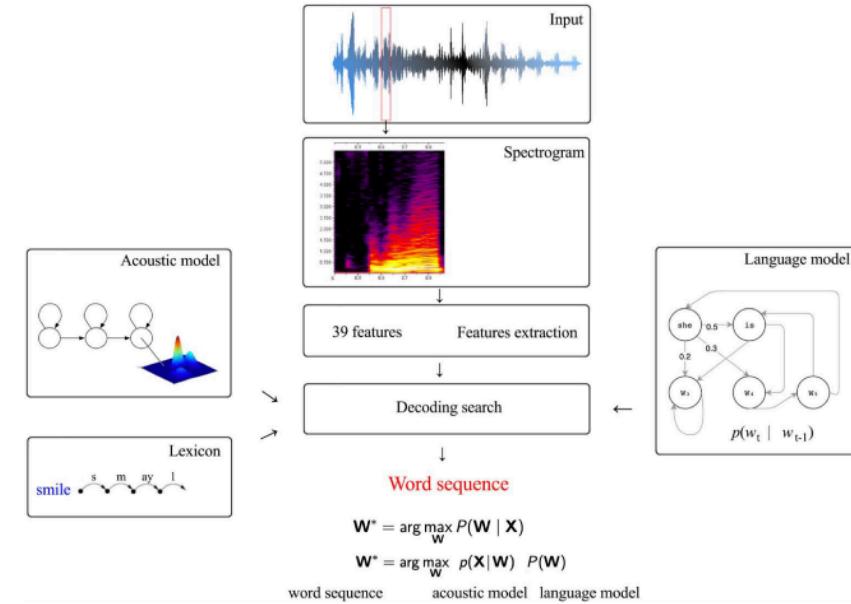
- $\sum_{i=1}^N \pi_i = 1$ $\sum_{j=1}^N a_{ij} = 1, \quad 1 \leq i \leq N$ $\sum_{k=1}^M b_j(k) = 1, \quad 1 \leq j \leq N$

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>

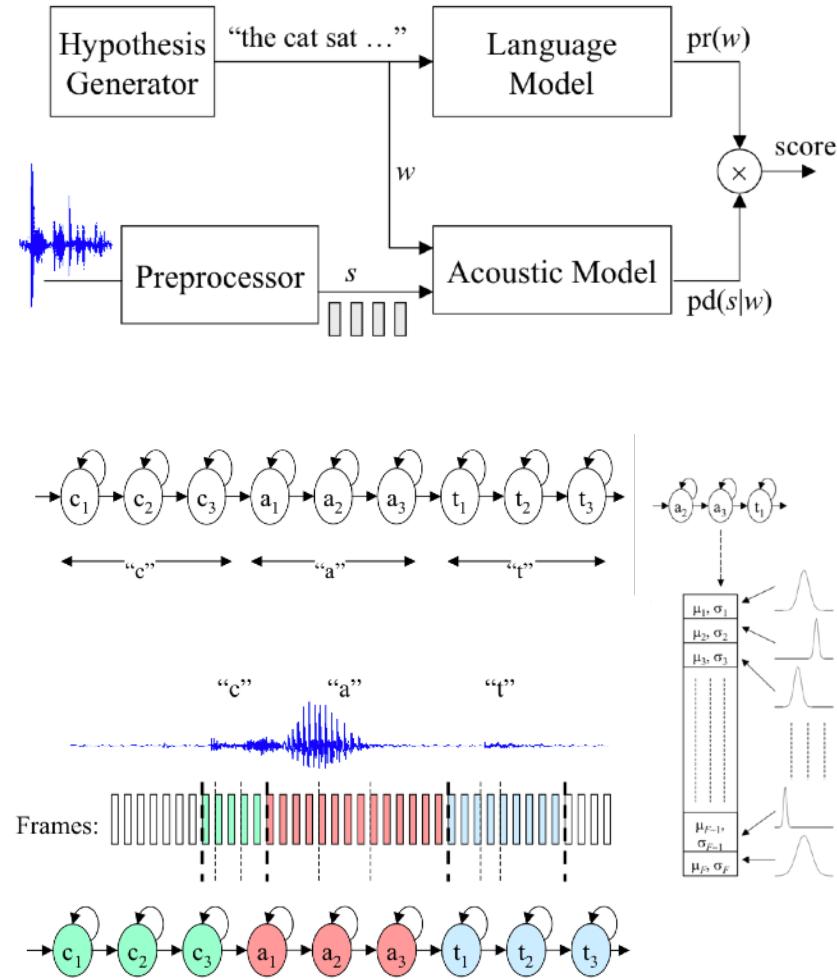
Automatic Speech Recognition (ASR) using GMM/HMM

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/ Δ/Δ^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

What you need to know

- What is a pitch
 - what is a temperament (tuning-system)
- 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 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 ?