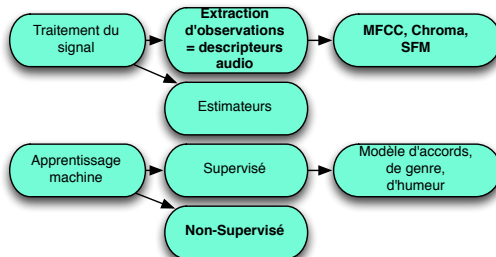
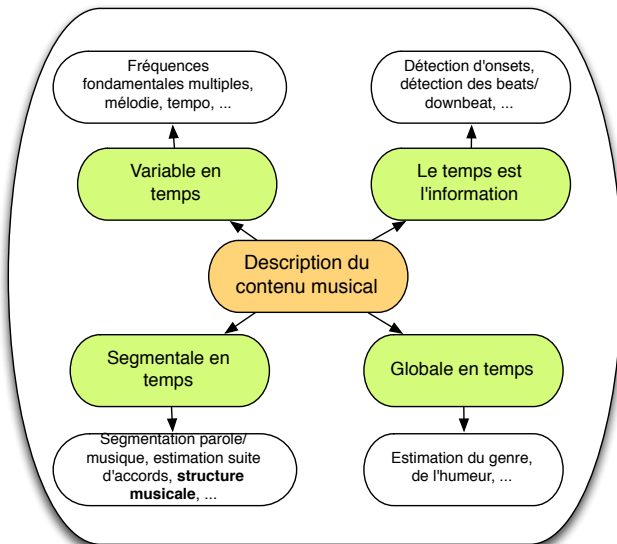


1- Introduction

Différents types de description du contenu musical



1- Introduction

Détection d'une structures musicale d'un morceau de musique

Objectifs

- Estimer une **structure** d'un morceau de musique
- Créer automatiquement un **résumé audio** représentatif du contenu du morceau

Applications

- Ecoute inter-active : création d'interface de lecture (player) interactif,
- Pré-écoute rapide d'un morceau
- **Exemples audio et vidéo**

Méthode

- Extraction de descripteurs audio
- Visualisation de la structure
- Estimation de la structure
 - Apprentissage non-supervisé (pas de pré-apprentissage possible)

The screenshot displays the 'INTERACTIVE PLAYER' interface. At the top, there's a search bar and a 'Musique' tab. The main player area shows the song 'Longtemps, longtemps (tu m'aimes en passant)' by Charléne Couture. Below the player is a 'RÉSULTATS (1463)' section with a table of results. The table has columns for 'Titre', 'Artiste', 'Album', and 'Durée'. The first row is 'Longtemps, longtemps (tu m'aimes en passant)' by Charléne Couture from the album 'Poèmes Rock' with a duration of 02:08. Other rows include 'Mister K.', 'Le Teneil d'Or', 'Last Night Thought', etc. To the right of the table are sections for 'Genres', 'Mouvements', 'Instrumentations', and 'TAG CLOUDS'. The 'TAG CLOUDS' section shows a list of instrumentations like 'Guitare électrique (1177)', 'Batterie Pop Legend Rock (1033)', etc.

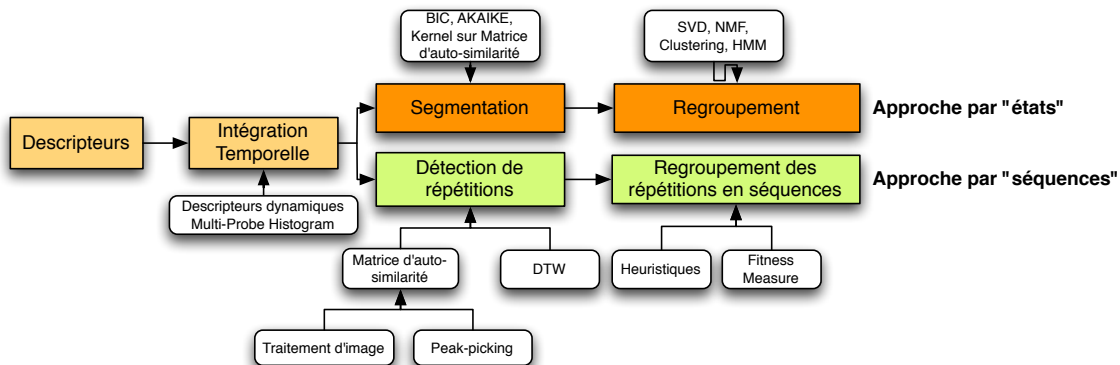
Titre	Artiste	Album	Durée
Longtemps, longtemps (tu m'aimes en passant)	Charléne Couture	Poèmes Rock	02:08
Mister K.	AMFON	Artificial Animals Riding On Neverland	02:07
Romantique - Pop/Rock - Guitare acoustique	AMFON	Artificial Animals Riding On Neverland	02:08
Le Teneil d'Or	AMFON	Artificial Animals Riding On Neverland	02:04
Last Night Thought	AMFON	Artificial Animals Riding On Neverland	02:04
Triste - Pop/Rock - Piano	AMFON	Artificial Animals Riding On Neverland	02:04
Let Me Put My Love into You	AC/DC	Back in Black	04:15
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:04
Blaise de Rio	AC/DC	Black Ice	03:07
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:02
Big Jack	AC/DC	Black Ice	03:07
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:02
Anything Goes	AC/DC	Black Ice	03:02
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	04:06
Swash in Grob	AC/DC	Black Ice	04:06
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:07
Wheels	AC/DC	Black Ice	03:08
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:03
Deafbe	AC/DC	Black Ice	03:03
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:03
Sherry May Day	AC/DC	Black Ice	03:10
Dynamique - Pop/Rock - Guitare électrique	AC/DC	Black Ice	03:10

source : Quacero project, MSSE-Orange interface

1- Introduction

Méthodes d'estimation de la structure

- 1) Extraction d'observations pertinentes du signal audio
 - **Descripteurs audio** : mise en évidence de différents contenus (timbre, harmonique, bruité, ...)
- 2) Analyse des observations afin de détecter une structure
 - Approche par **états**
 - **segmentation** temporelle et
 - **regroupement** des segments homogènes identiques
 - Approche par **séquences**
 - **détection des répétitions** non-homogènes et
 - regroupement des segments répétés en séquences



2- Descripteurs audio

Les descripteurs audio

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

- Valeurs numériques extraites du signal audio dont le but est de représenter une propriété particulière de son contenu
 - Tout est dans la forme d'onde, dans la TFCT, difficile à lire, trop grande dimension
- Contrainte :
 - on veut le même nombre de dimensions pour toutes les données
- Extraction ?
 - Algorithme d'estimation
 - Opérateurs mathématique

2- Descripteurs audio

Introduction

Les descripteurs audio

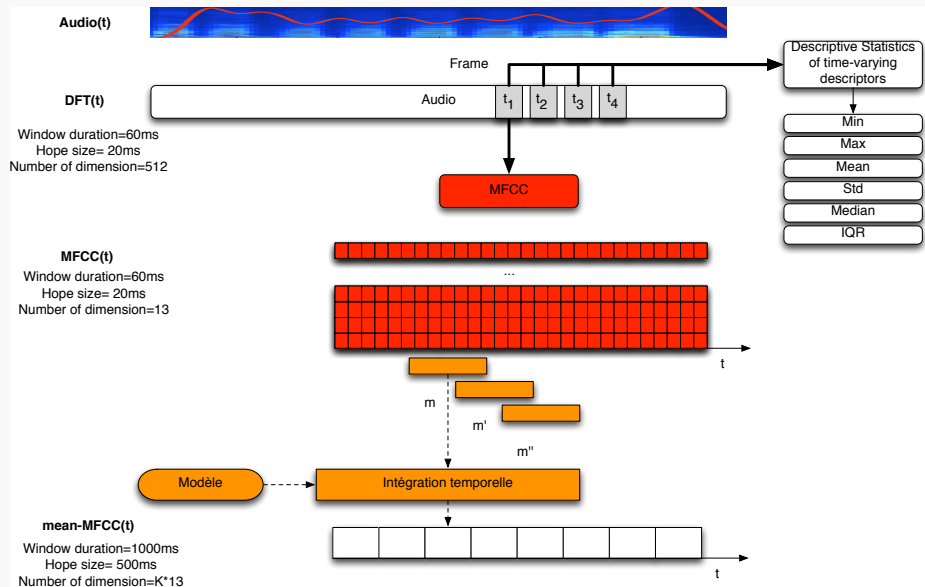
[G. Peeters. A large set of audio features for sound description (similarity and classification) in the cuidado project. Cuidado project report, Ircam, 2004.]

- Différentes **formes** :
 - **scalaire** : Centroïde spectral, étendue spectrale, fréquence fondamentale, spectral roll-off, spectral flux, zero-crossing rate, RMS, ...
 - **vecteur** : Mel Frequency Cepstral Coefficients, coefficients LPC, coefficients PLP ...
- Différentes **temporalité** :
 - représente une **trame** du signal audio → descripteurs "instantanés"
 - représente le résumé du contenu d'un **ensemble local de trame** → texture windows
 - représente **globalement** le signal audio
- Mise en évidence de différents **contenus** (, harmonique, bruité, ...)
 - contenu **timbral** : Mel Frequency Cepstral Coefficients, coefficients LPC, coefficients PLP ...
 - contenu **harmonique** : Pitch Class Profiles/ Chroma ...
 - contenu **bruité** : Spectral Flatness Measure
 - contenu **rythmique** : ...

2- Descripteurs audio

Introduction

Les descripteurs audio

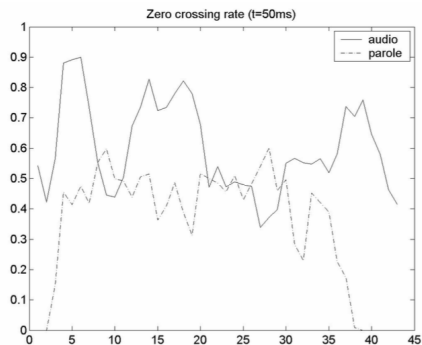
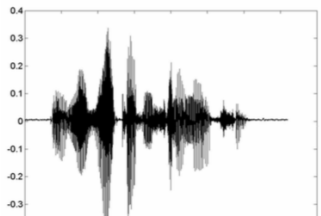
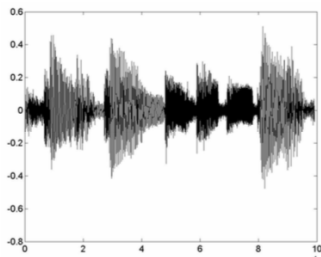


2- Descripteurs audio

Taux de passage par zéro

Taux de passage par zéro / zero-crossing rate (zcr)

- Mesure le nombre de fois que la forme d'onde croise l'axe zéro
 - $zcr = 0.5 \sum_{n=1}^N |sign(x(n)) - sign(x(n-1))|$
- Utilisation :
 - permet de distinguer les signaux bruités \rightarrow zcr élevé
 - permet de distinguer les signaux harmoniques \rightarrow zcr bas



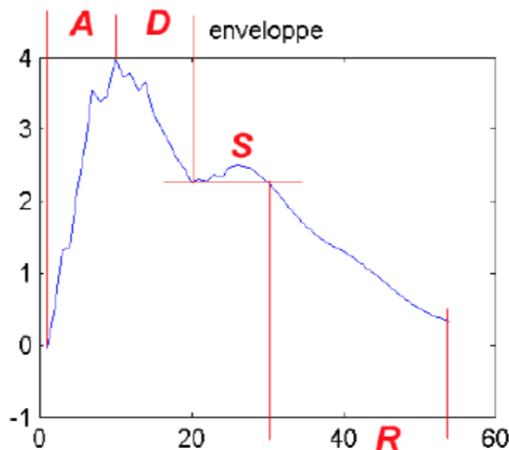
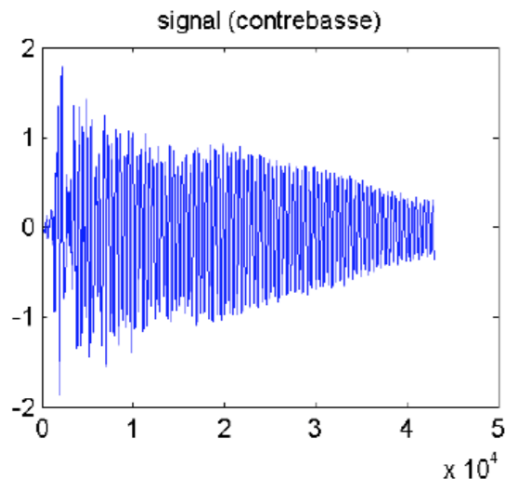
source : Gaël Richard

2- Descripteurs audio

Enveloppe ADSR

Enveloppe ADSR (Attack, Decay, Sustain, Release)

- Modèle représentant l'évolution (l'enveloppe) d'énergie d'une note de musique
- Utilisation :
 - permet de distinguer les attaques rapides (sons percussifs) / lentes
 - permet de distinguer les décroissances rapides (sons non-tenus) / lentes (sons tenus)



2- Descripteurs audio

Description du spectre (barycentre, étendue spectral)

Description du spectre (barycentre, étendue spectral)

- **Centroid spectral**

- $cs = \frac{\sum_k f_k A_k}{\sum_k A_k}$

- Utilisation :

- permet de distinguer les sons terne des sons brillant

- **Etendue spectral**

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

- Utilisation :

- permet de distinguer les sons pauvres des sons riches

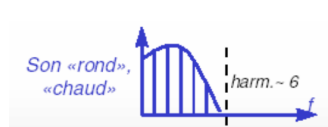
- **Flux spectral**

- Mesure la variation temporel du spectre

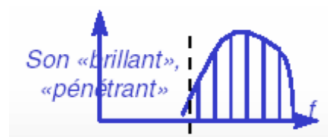
- $fs = \sum_k (A_k(t) - A_k(t-1))^2$

- Utilisation :

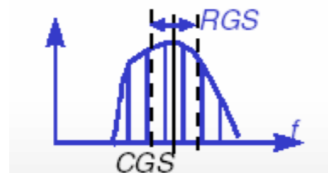
- permet de distinguer les sons pauvres des sons riches



source : Gaël Richard



source : Gaël Richard



source : Gaël Richard

2- Descripteurs audio

Mel Frequency Cepstral Coefficients (MFCCs)

Objectif

- décrire la forme du spectre (du timbre) d'un signal à l'aide d'un nombre réduit de coefficients

Cepstre complexe

- Cepstre complexe** $c(\tau)$:

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

- τ est appelé "céfrence"
- $x(t) \xrightarrow{TF} X(\omega) \xrightarrow{\log} \log(X(\omega)) \xrightarrow{TF^{-1}} c(\tau)$

2- Descripteurs audio

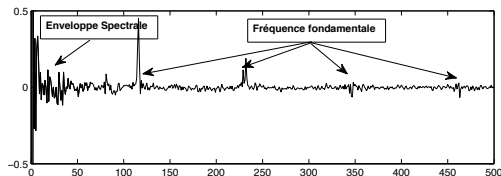
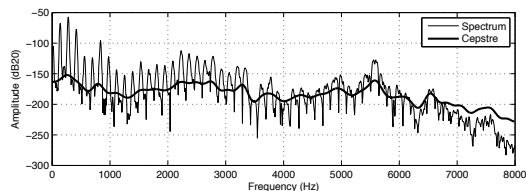
Mel Frequency Cepstral Coefficients (MFCCs)

Cepstre complexe

- Modèle source/ filtre :
 - Source : signal périodique
 - Filtre : résonant/ anti-résonant

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

$$\xrightarrow{TF} X(\omega) = E(\omega) \cdot G(\omega) \quad (2)$$



$$\xrightarrow{\log} \log(X(\omega)) = \underbrace{\log(E(\omega))}_{\text{variation rapide à travers } \omega} + \underbrace{\log(G(\omega))}_{\text{variation lente à travers } \omega} \quad (3)$$

$$\xrightarrow{TF^{-1}} TF^{-1} [\log(X(\omega))] = \underbrace{TF^{-1} [\log(E(\omega))]}_{\text{énergie aux céfrenes } \tau \gg} + \underbrace{TF^{-1} [\log(G(\omega))]}_{\text{énergie aux céfrenes } \tau \ll}$$

2- Descripteurs audio

Mel Frequency Cepstral Coefficients (MFCCs)

Cepstre réel

- **Cepstre réel** :
 - Cepstre calculé sur la partie réelle du log-spectrum

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

$$\log(X(\omega)) = \log(A(\omega)) + j\phi(\omega) \quad (4)$$

$$\Re(\log(X(\omega))) = \log(A(\omega))$$

$$\begin{aligned} \text{cepstre réel} &= TF^{-1} [\Re(\log(X(\omega)))] \\ &= TF^{-1} [\log(A(\omega))] \end{aligned} \quad (5)$$

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

- Le spectre d'amplitude étant réel et symétrique
 - sa TF se réduit à sa partie réelle
 - donc à la projection de $\log(A(\omega))$ sur un ensemble de cosinus \rightarrow DCT

2- Descripteurs audio

Mel Frequency Cepstral Coefficients (MFCCs)

Mel Frequency Cepstral Coefficients (MFCCs)

- **Mel Frequency Cepstral Coefficient :**
 - Cepstre réel calculé sur un spectre d'énergie exprimé en convertissant l'énergie $|X(\omega)|^2$ en échelle perceptive (échelle de Mel)
- Pourquoi ?
 - La transformée de Fourier :
 - décomposition sur une série de sinusoides linéairement espacées (10Hz, 20Hz, 30Hz, ... Hz)
 - L'oreille :
 - décomposition sur une série de filtres de fréquences logarithmiquement espacé (10, 20, 40, 80, ... Hz).
 - meilleure résolution en basses fréquences que en hautes fréquences.
 - résonances de l'enveloppe spectrale sont plus rapprochées en basse fréquence.
 - MFCCs permet une représentation plus compacte que le cepstre réel
- Comment ?
 - On utilise des échelles dites perceptives : échelles de Mel, de Bark, filtres ERB, Gamma tone
- Utilisation ?
 - Les coefficients les plus utilisés dans le monde de la reconnaissance audio : parole, musique, sons environnementaux, ...

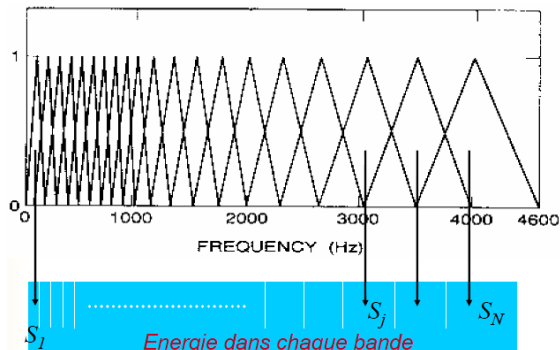
2- Descripteurs audio Mel Frequency Cepstral Coefficients (MFCCs)

Mel Frequency Cepstral Coefficients (MFCCs)

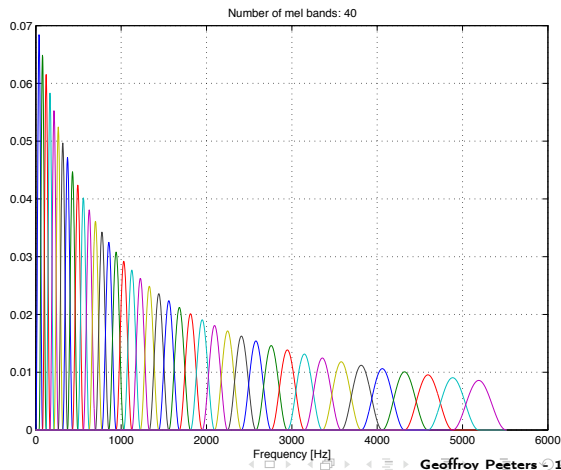
- Echelle de Mel :

$$M = f \text{ pour } f < 1000\text{Hz}$$

$$M = f_c \left(1 + \log_{10} \left(\frac{f}{f_c} \right) \right) \text{ pour } f \geq 1000\text{Hz} \quad (6)$$



source : Gaël Richard

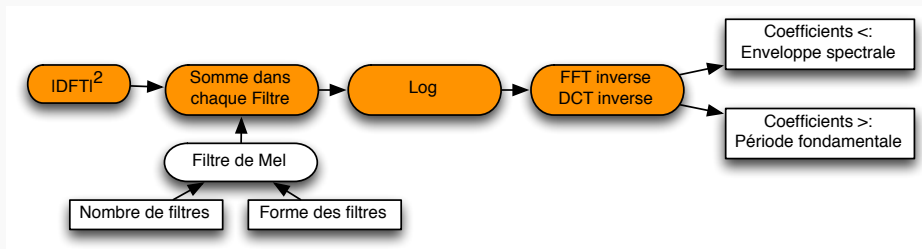


2- Descripteurs audio

Mel Frequency Cepstral Coefficients (MFCCs)

Mel Frequency Cepstral Coefficients (MFCCs)

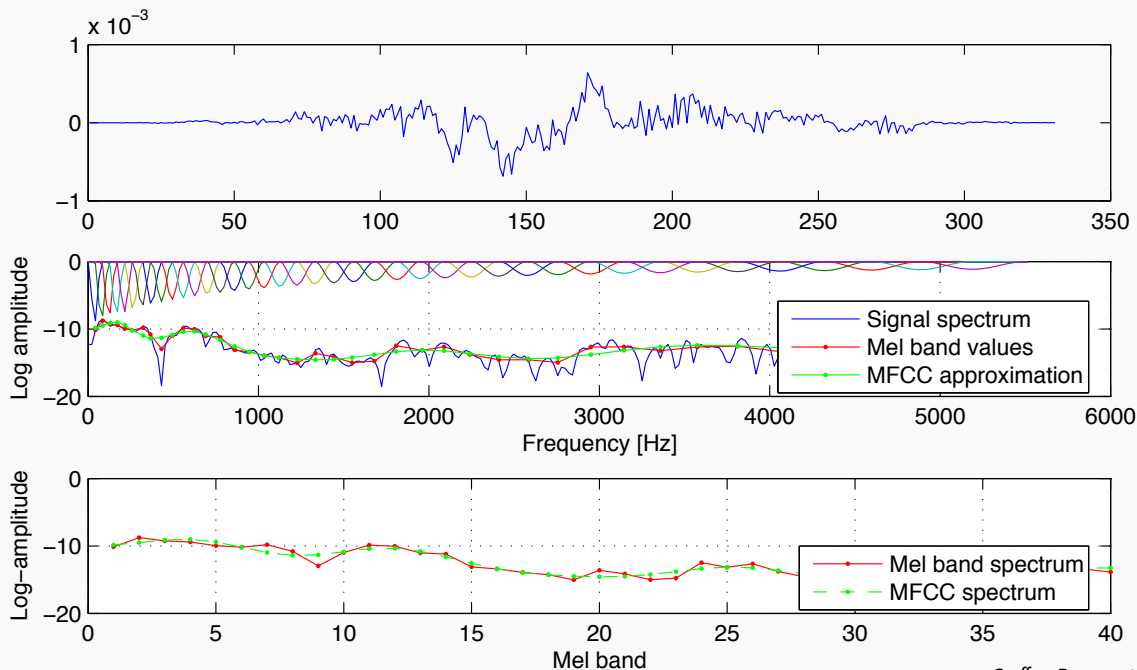
- Calcul du spectre de puissance : $|X(\omega)|^2$
- Calcul des filtres de Mel : $H_b(\omega)$ avec $b \in [1, B]$
 - choix du nombre de filtres B : 40
 - choix de la forme des filtres : triangulaire, hanning, tanh, ...
- Conversion du spectre de puissance en bandes de Mel : $S(b) = \sum_{\omega} |X(\omega)|^2 \cdot H_b(\omega)$
- Passage en échelle logarithmique : $\log(S(b))$
- Calcul de la IFFT (ou de la IDCT) :
- Sélection des coefficients de la IDCT proches de zéro (jusqu'à 13)
 - les coefficients proches de zéro représentent la décomposition du spectre en échelle de Mel sur un ensemble de cosinus à variation lente



2- Descripteurs audio

Mel Frequency Cepstral Coefficients (MFCCs)

Exemple de calcul de MFCCs



2- Descripteurs audio

Chroma - Pitch Class Profile (PCP)

Définition des Chroma - Pitch Class Profile (PCP)

- **Objectif :**

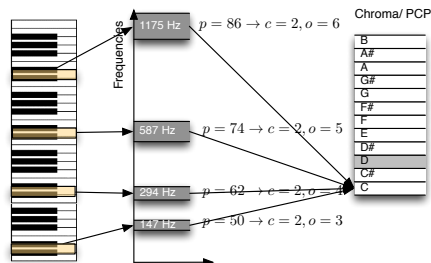
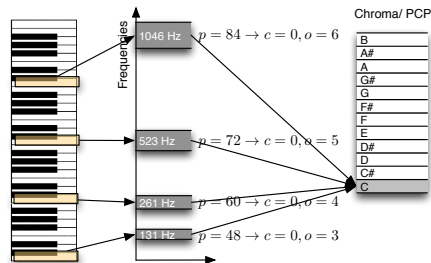
- le spectre à l'instant n : $X(k, n)$
- représenter son contenu harmonique sous forme d'un vecteur : $C(c, n)$ $c \in [0, 12[$

- Utilisations :

- reconnaissance de tonalité,
- reconnaissance de suite d'accords,
- détection de "cover versions"

- Shepard-1964 :

- représenter la hauteur d'une note p comme une structure bi-dimensionnelles :
- $p = c + o \cdot 12$
 - le chroma c (classe de hauteur).
 - la hauteur tonale o (numéro d'octave),



2- Descripteurs audio

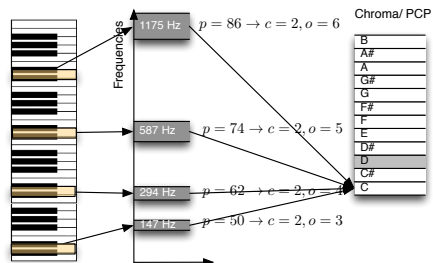
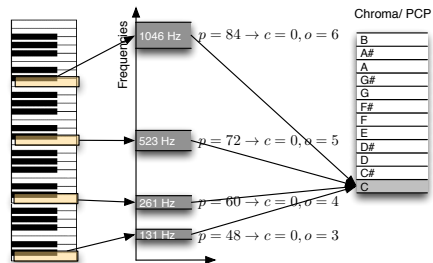
Chroma - Pitch Class Profile (PCP)

Calcul des Chromas - Pitch Class Profile (PCP)

- Relation entre les fréquences f_k de la DFT et les hauteurs de note p (hauteurs de demi-tons en échelle de notes MIDI)

- $p(f_k) = 12 \log_2 \left(\frac{f_k}{440} \right) + 69, p \in \mathbb{R}^+$
- $f(p) = 440 \cdot 2^{\frac{p-69}{12}}$

- Calcul des chromas $C(c, n)$
 - On additionne toutes les valeurs du spectre $X(k, n)$ tel que f_k correspondent à un c donné
 - Hard-mapping
 - Soft-mapping



2- Descripteurs audio

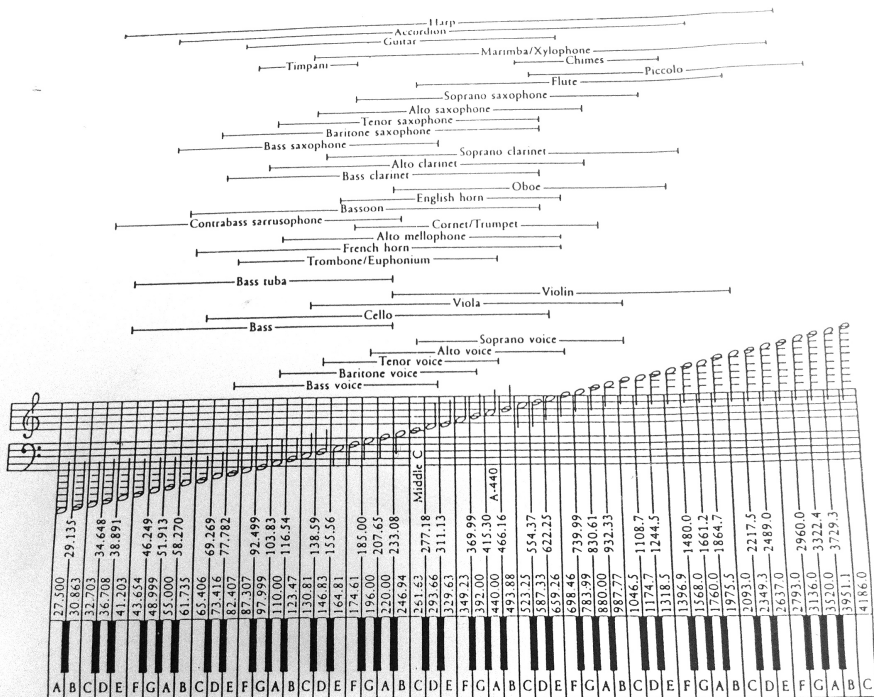
Chroma - Pitch Class Profile (PCP)

Calcul des Chromas - Pitch Class Profile (PCP)

- Résolution fréquentielle ?
 - Elle doit permettre la séparation des notes voisines
 - On définit la largeur (à -6 dB) : $Bw = \frac{Cw}{L_{sec}}$
 - Si f_{\min} (la fréquence la plus basse considérée dans le secteur) est 50 Hz
 - on veut séparer G#1 (51.91Hz) et A1 (55Hz) $\rightarrow L_{sec} = \frac{Cw}{Bw} = \frac{2.35}{3.0869Hz} = 0.7613s$
 - Si f_{\min} est 100 Hz
 - on veut séparer G#2 (103.82Hz) de A2 (110Hz) $\rightarrow L_{sec} = \frac{Cw}{Bw} = \frac{2.35}{6.1738Hz} = 0.3806s$
- Deux possibilités :
 - Choisir L_{sec} en fonction f_{\min}
 - Choisir f_{\min} en fonction de L_{sec}

2- Descripteurs audio

Chroma - Pitch Class Profile (PCP)

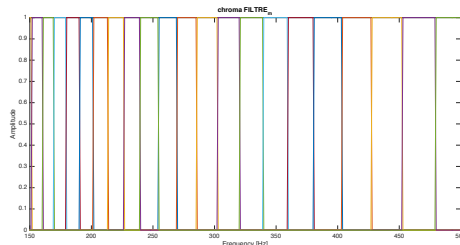


2- Descripteurs audio

Chroma - Pitch Class Profile (PCP)

Hard-mapping

- Hard-mapping ?
 - Une fréquence f_k de la DFT contribue uniquement à la note la plus proche
 - Par exemple,
 - l'énergie à $f_k=452$ Hz ($p(f_k)=69.4658$) contribue entièrement à la note $p=69$ ($c=10$)
 - alors que $f_k=453$ Hz ($p(f_k)=69.5041$) à $p=70$ ($c=11$).
- Création d'un banc de filtres $H_{p'}$ centrés sur les hauteurs de demi-tons $p' \in [43, 44, \dots, 95]$:



2- Descripteurs audio

Chroma - Pitch Class Profile (PCP)

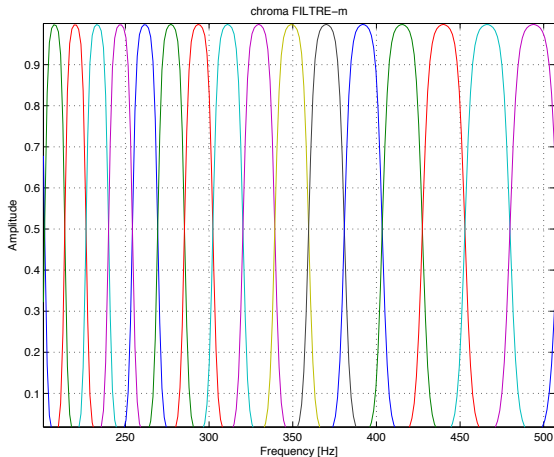
Soft-mapping

- Création d'un banc de filtres $H_{p'}$ centrés sur les hauteurs de demi-tons $p' \in [43, 44, \dots, 95]$:
 - Chaque filtre est défini par la fonction

$$H_{p'}(f_k) = \frac{1}{2} \tanh(\pi(1 - 2x)) + \frac{1}{2}$$

dans lequel $x =$ distance relative entre centre du filtre et fréquences de la TF
 $x = R |p' - p(f_k)|$.

- Les filtres sont équi-répartis et symétriques sur l'échelle logarithmique des hauteurs de demi-tons, non-nulles entre $p' - 1$ et $p' + 1$ et à valeur maximale en p' .



2- Descripteurs audio

Chroma - Pitch Class Profile (PCP)

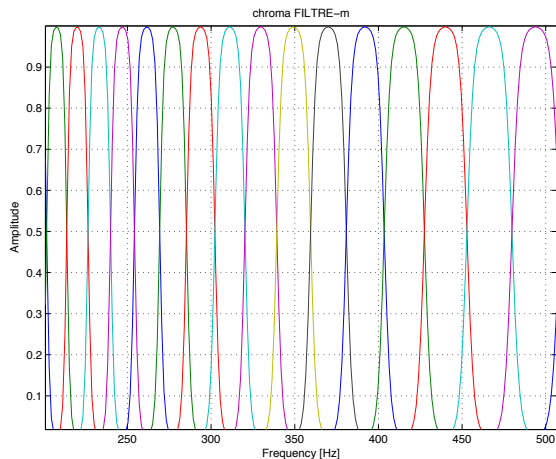
Calcul des Chromas - Pitch Class Profile (PCP)

- La valeur du spectre de hauteur de demi-ton $N(n')$ est obtenue en multipliant les valeurs de la transformée de Fourier $A(f_k)$ par l'ensemble des filtres $H_{n'}$:

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

- Le mapping entre les hauteurs de demi-tons n et les classes de hauteurs de demi-tons (chroma) c est défini par $c(p) = \text{mod}(p, 12)$.
- La valeur du vecteur de chroma est obtenue en additionnant les valeurs de classes de hauteur équivalentes

$$C(c) = \sum_{p' \text{ tel que } c(p')=l} P(n')$$



2- Descripteurs audio

Chroma - Pitch Class Profile (PCP)

Limitations des Chromas - Pitch Class Profile (PCP)

- Présence des harmoniques supérieures de chaque note
 - En pratique pour une note C on a pas $[1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$
 - mais plutôt $[a_1 + a_2 + a_4, 0, 0, 0, a_5, 0, 0, a_4, 0, 0, 0, 0]$
- Influence de l'enveloppe spectrale

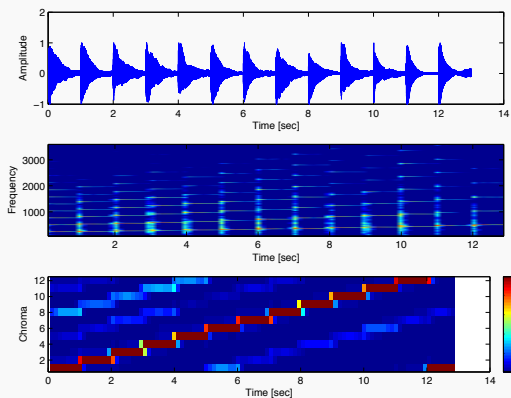
Pitch	Harmonic	Frequency f_μ	MIDI-scale m_μ	Chroma/PCP p
c3	f_0	130.81	48	1 (=c)
	$2f_0$	261.62	60	1 (=c)
	$3f_0$	392.43	67.01	8.01 (\simeq g)
	$4f_0$	523.25	72	1 (=c)
	$5f_0$	654.06	75.86	4.86 (\simeq e)

2- Descripteurs audio

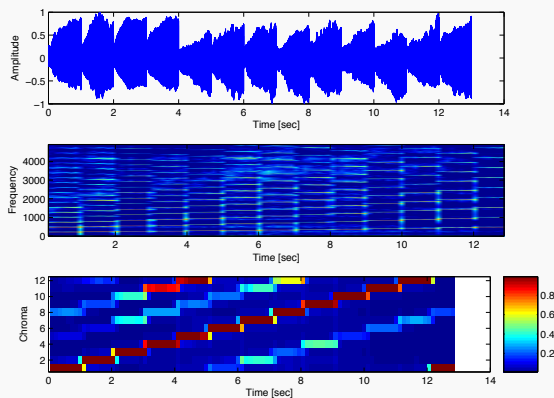
Chroma - Pitch Class Profile (PCP)

Limitations des Chromas - Pitch Class Profile (PCP)

Exemple piano



Exemple violon



Représentation visuelle de la structure temporelle de la musique

3- Représentation visuelle de la structure temporelle de la musique

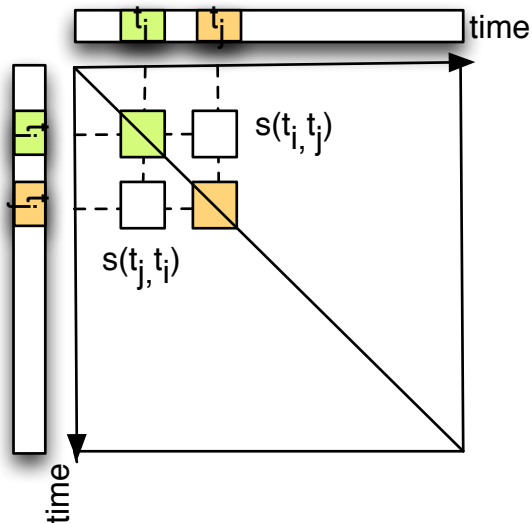
La matrice d'auto-similarité

La matrice d'auto-similarité

- Similarité entre deux instants t_i et t_j
- Similarité entre les observations du signal aux trames i et j : $s(t_i, t_j) = s(\underline{d}^i, \underline{d}^j)$
- Matrice d'auto-similarité = les valeurs $s(t_i, t_j)$ sont représentées sous forme d'une matrice $\underline{S} = s(t_i, t_j) \quad \forall i, j$

Lecture/ interprétation

- Une valeur élevée dans $S(t_i, t_j)$ représente une similarité importante entre les instants t_i et t_j .
- Si $t_i \simeq t_{i+1} \simeq t_{i+2}$ nous observons un **bloque homogène**
- Si une **séquence de temps** $t_i, t_{i+1}, t_{i+2}, \dots$ est similaire à une séquence de temps $t_j, t_{j+1}, t_{j+2}, \dots$ nous observons une diagonale supérieure/ inférieure dans S .

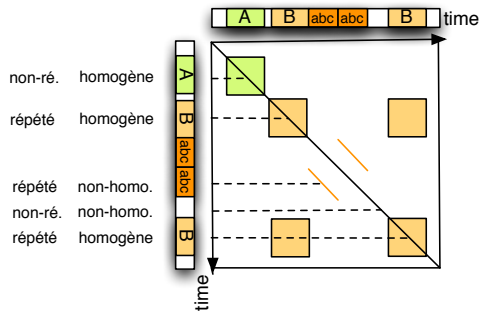


3- Représentation visuelle de la structure temporelle de la musique

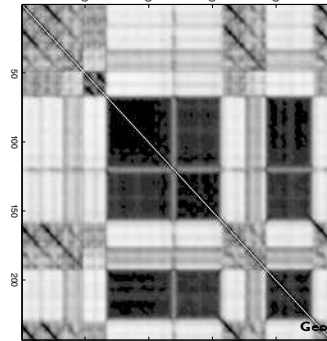
Hypothèses concernant la macro-structure d'un morceau

Hypothèse 1 : homogénéité

- Hypothèse : le morceau est formé d'une succession de segments temporels **homogènes** $t_i \simeq t_{i+1} \simeq t_{i+2}, \dots$ et de segments non homogènes
 - homogène? contenant une information similaire au sens d'un critère d'observation)
 - "A" et "B" sur la Figure
- Exemple : arrangements d'un couplet ou d'un refrain
- Méthode : **approche par "état"**



	homogène	non-homogène
répété	état	séquence
non-répété	état	état

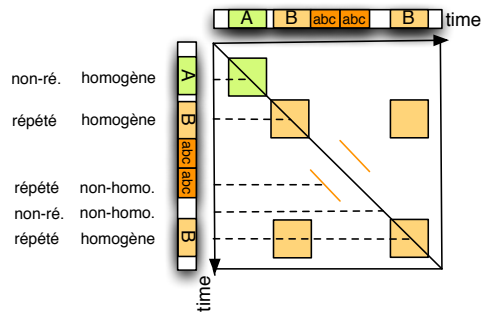


3- Représentation visuelle de la structure temporelle de la musique

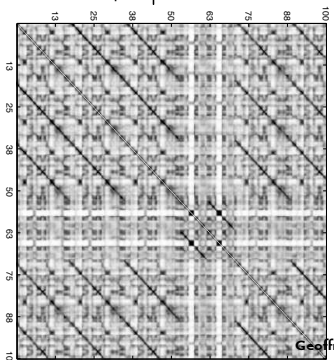
Hypothèses concernant la macro-structure d'un morceau

Hypothèse 2 : répétition

- Hypothèse : le morceau renferme des **répétitions** temporelles.
 - elles peuvent correspondre à des répétitions de segments **homogènes**
 - $\{t_j, t_{j+1}, t_{j+2}\} \simeq \{t_i, t_{i+1}, t_{i+2}\}$ et $t_i \simeq t_{i+1} \simeq t_{i+2}$
 - "B" dans la figure
 - Méthode : **approche par "état"**
 - elles peuvent correspondre à des répétitions de segments **non homogènes**
 - $\{t_j, t_{j+1}, t_{j+2}\} \simeq \{t_i, t_{i+1}, t_{i+2}\}$ et $t_i \neq t_{i+1} \neq t_{i+2}$
 - séquence "abc" dans la Figure
 - Méthode : **approche par "séquence"**



	homogène	non-homogène
répété	état	séquence
non-répété	état	état



3- Représentation visuelle de la structure temporelle de la musique

3- Représentation visuelle de la structure temporelle de la musique

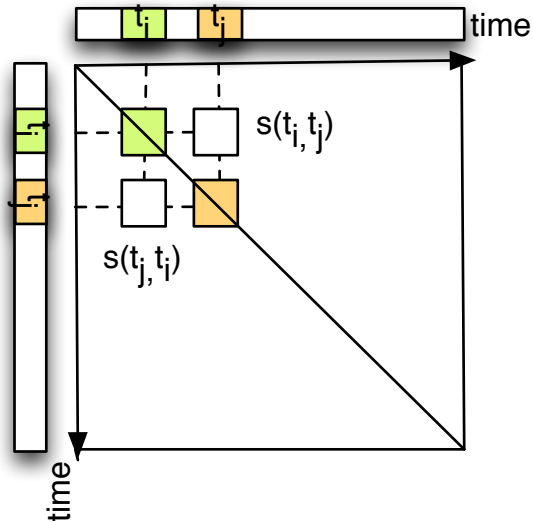
Matrice d'auto-similarité/distance (temps, temps)

Matrice d'auto-similarité/distance (temps, temps)

- Similarité entre deux instants t_i et t_j
- Similarité entre les observations du signal à deux trames i et j : $s(t_i, t_j) = s(\underline{d}^i, \underline{d}^j)$
- Descripteurs audio multi-dimensionnels $\underline{d} = d_k \quad k \in K$

Choix d'une distance

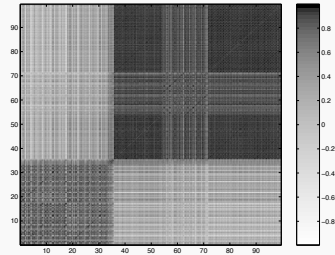
- Distance euclidéenne : $\sqrt{\sum_k (d_k^i - d_k^j)^2}$
- Corrélation : $\sum_k (d_k^i \cdot d_k^j)$
- Distance cosinusoidale : $\frac{\sum_k (d_k^i \cdot d_k^j)}{\sqrt{\sum_k (d_k^i)^2} \sqrt{\sum_k (d_k^j)^2}}$
- Correlation Pearson : $\frac{\sum_k (d_k^i - \mu^i) \cdot (d_k^j - \mu^j)}{\sqrt{\sum_k (d_k^i - \mu^i)^2} \sqrt{\sum_k (d_k^j - \mu^j)^2}}$
- ...



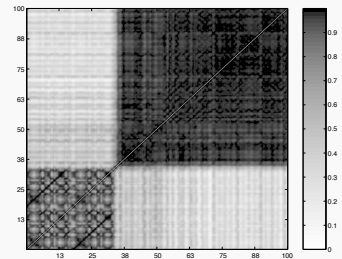
3- Représentation visuelle de la structure temporelle de la musique

Matrice d'auto-similarité/distance (temps, temps)

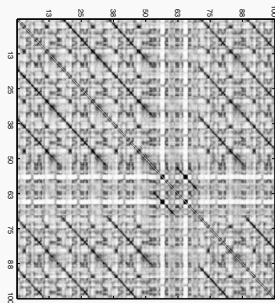
MFCC



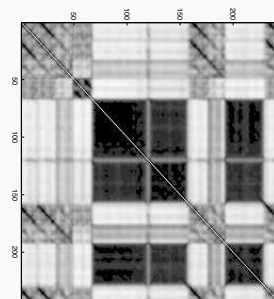
Modulation spectrum 1



Chroma



Modulation spectrum 2

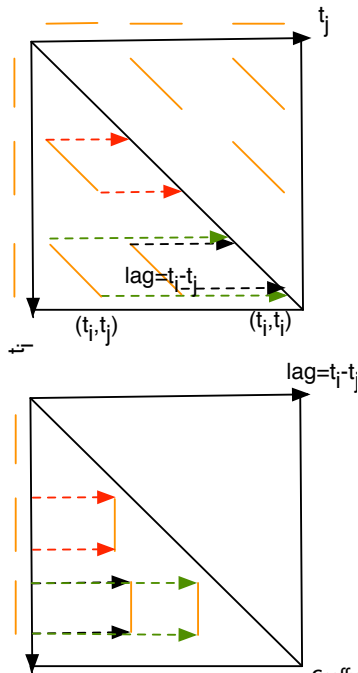


3- Représentation visuelle de la structure temporelle de la musique

Matrice d'auto-similarité/distance (temps,lag)

Matrice d'auto-similarité/distance (temps,lag)

- Une valeur élevée dans $S(t_i, t_j)$ représente une similarité importante entre les instants t_i et t_j .
- Si une séquence de temps $t_i, t_{i+1}, t_{i+2}, \dots$ est similaire à une séquence de temps $t_j, t_{j+1}, t_{j+2}, \dots$ nous observons une diagonale supérieure/ inférieure dans S .
- **Lag** = distance entre la répétition (démarrant en t_i) et la séquence originale (démarrant en t_j)
 - cette distance est donnée par la projection de t_i sur la diagonale principale de la matrice : $t_i - t_j$
 - souvent constante
- Matrice de lag :
 $L = L(t_i, lag_{ij}) = S(t_i, t_i - t_j)$
 - les diagonales dans une matrice (temps, temps)
 - deviennent des lignes verticales dans une matrice (temps, lag)

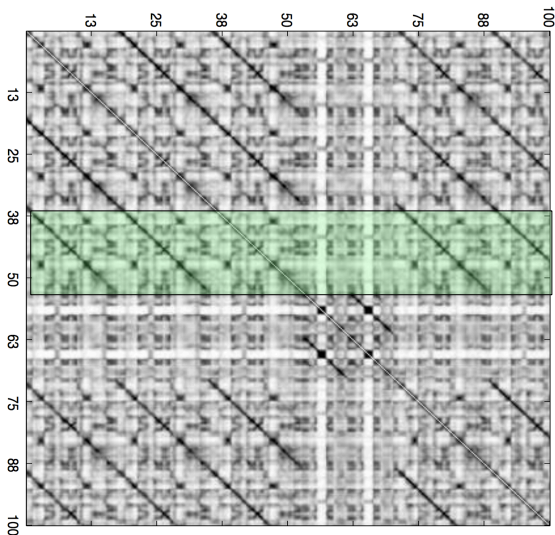


3- Représentation visuelle de la structure temporelle de la musique Génération de résumé audio par méthode du "summary score"

Génération de résumé audio par méthode du "summary score"

[M. Cooper and J. Foote. Automatic music summarization via similarity analysis. In Proc. of ISMIR, Paris, France, 2002.]

- Recherche du segment temporel continu représentant au mieux le contenu d'un morceau de musique selon un critère de similarité → création de "previews" musicaux

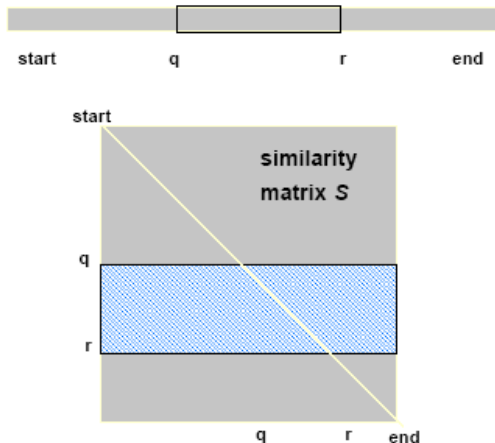


3- Représentation visuelle de la structure temporelle de la musique

Génération de résumé audio par méthode du "summary score"

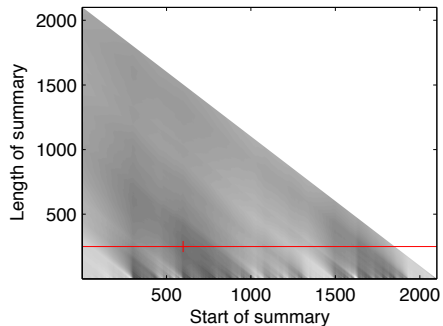
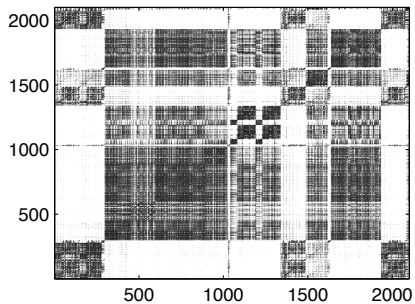
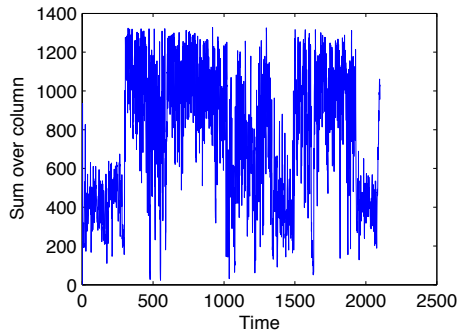
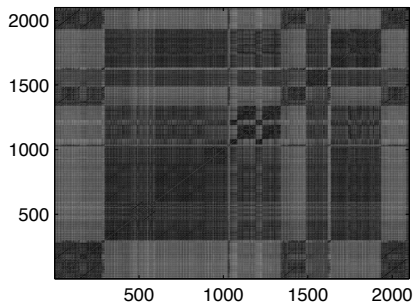
Génération de résumé audio par méthode du "summary score"

- Recherche du segment démarrant en q de durée $L = r - q$ expliquant le maximum de répétitions
- Similarité moyenne **de l'instant** q avec tous les temps du morceau
 - $\frac{1}{N} \sum_{n=1}^N S(q, n)$
- Similarité moyenne **du segment** $[q, r]$ (de longueur $L = r - q$) avec tous les temps du morceau
 - $s(q, L) = \frac{1}{LN} \sum_{m=q}^r \sum_{n=1}^N S(m, n)$
- Pour un L donné, nous cherchons q maximisant $s(q, L)$
 - $q_L = \operatorname{argmax}_{1 \leq i \leq N-L} s(i, L)$
- Variante : pour favoriser la détection de résumés en début de morceau,
 - ajout d'une pondération $w(n)$ fonction décroissante du temps
 - $s(q, L) = \frac{1}{LN} \sum_{m=q}^r \sum_{n=1}^N w(n) S(m, n)$



source : [Cooper and Foote, 2002, ISMIR]

3- Représentation visuelle de la structure temporelle de la musique Génération de résumé audio par méthode du "summary score"



4- Segmentation temporelle d'un flux de descripteurs

4- Segmentation temporelle d'un flux de descripteurs

Critère BIC (Bayes Information Criteria)

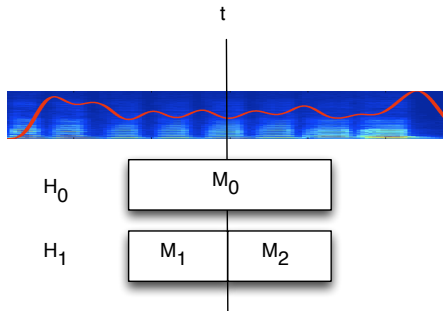
Critère BIC (Bayes Information Criteria)

- Pour chaque temps t (potentiellement instant de rupture) on compare deux hypothèses
 - H_0 : le signal obéit au même modèle probabiliste de part et d'autre de t , modèle noté $M_0(\mu_0, \Sigma_0)$
 - H_1 : il y a un changement de modèle en t , deux modèles différents $M_1(\mu_1, \Sigma_1)$ et $M_2(\mu_2, \Sigma_2)$
- Critère Delta BIC

$$\Delta BIC = R(t) - \lambda P$$

$$R(t) = \frac{1}{2}(N \log(|\Sigma_0|) - t \log(|\Sigma_1|) - (N - t) \log(|\Sigma_2|))$$

- si $\Delta BIC > 0$, H_1 est vérifiée
- paramètres :
 - P : proportionnel à la différence des nombres de paramètres estimés pour chaque hypothèse
 - λ : facteur de pénalité choisi tel que $\Delta BIC > 0$ si H_1 est vérifiée



4- Segmentation temporelle d'un flux de descripteurs

Convolution de la matrice d'auto-similarité par un noyau en damier

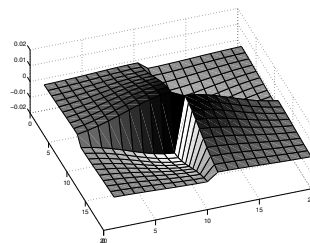
Convolution de la matrice d'auto-similarité par un noyau en damier

[J. Foote. *Automatic audio segmentation using a measure of audio novelty*. In *Proc. of IEEE ICME, New York City, NY, USA, 2000.*]

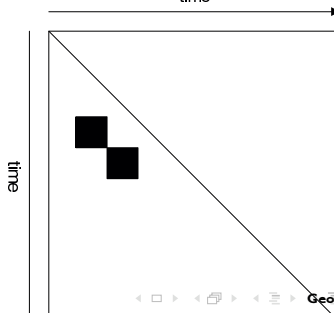
- Méthode "Novelty Curve" [Foote, 2000, ICME]
- Approche plus robuste
- Convolution de la matrice de similarité \underline{S} par un noyau prenant en compte
 - la similarité inter-segment (homogénéité) et
 - la dis-similarité entre **segments** gauches et droites
 - "checker-board" kernel :

$$C = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

- La valeur de la diagonale de la matrice "filtrée" mesure la similarité/ dis-similarité des **régions** gauches et droites



source : [Foote, 2000, ICME]
time



4- Segmentation temporelle d'un flux de descripteurs

Convolution de la matrice d'auto-similarité par un noyau en damier

Exemple

