

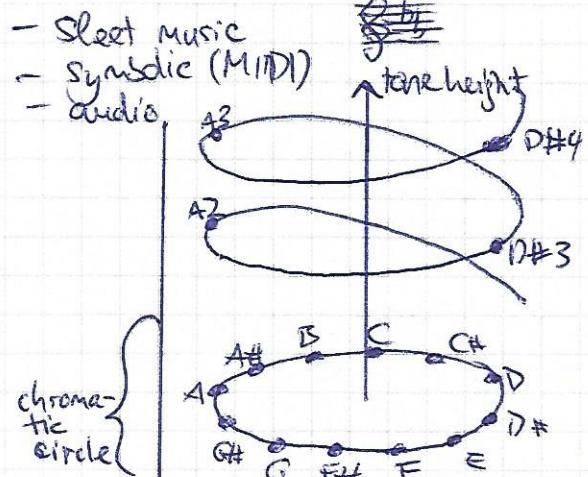


Exam preparation

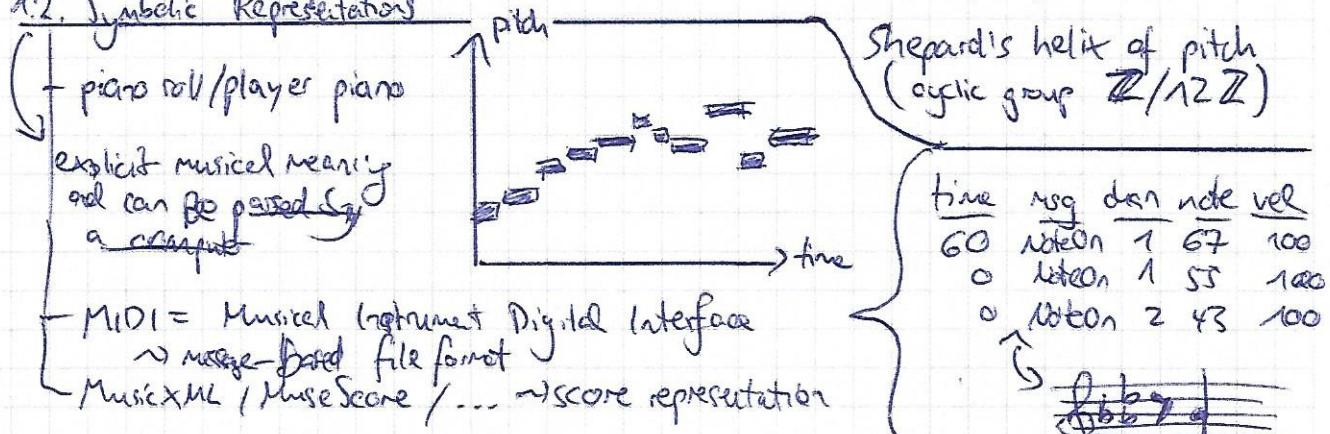
Chapter 1. Music Representations

1.1. Sheet Music Representation

- pitches map into a pitch class
- an octave doubles or halves the frequency
- here: twelve-tone equal-tempered scale
- Scientific Pitch Notation:
 $A4 \approx 440 \text{ Hz}$
- chroma = "color" of a pitch, pitch class



1.2. Symbolic Representations



1.3. Audio Representation

music as acoustic waves (waveforms)

- period = cycle length
- frequency $\geq \frac{1}{\text{period}}$ (Hz)
- amplitude = peak mean deviation
- phase = x s.t. $\sin(x) = 0$
- audible frequency range: 20 Hz - 20 kHz
- equal-tempered scale:

$$f(p) = 2^{\frac{p-69}{12}} \cdot 440 \text{ Hz}$$

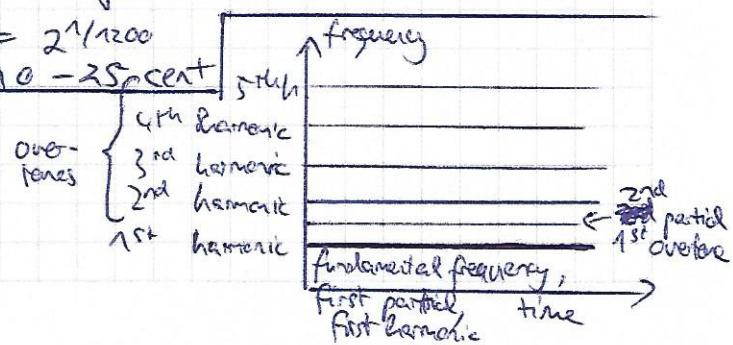
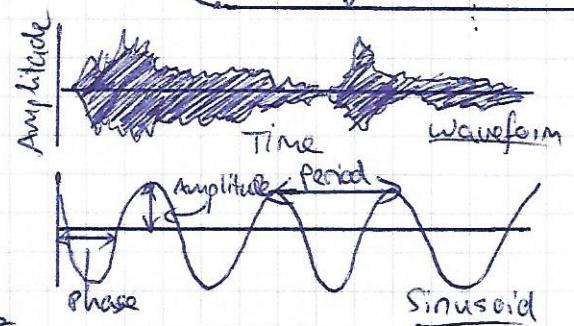
$$\text{in particular } f(p+12) = 2f(p)$$

$$\frac{f(p+1)}{f(p)} = 2^{1/12} = \sqrt[12]{2} \approx 1.059 = 100 \text{ cent}$$

$$\Leftrightarrow \text{i.e. } 1 \text{ cent} = 2^{1/1200}$$

\Rightarrow just noticeable difference: $10 - 25 \text{ cent}$

- partials (overtones/harmonics)
are closely related to timbre



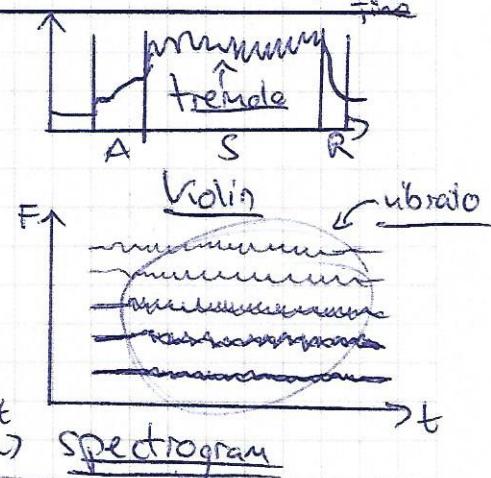
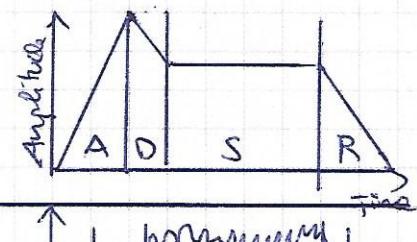
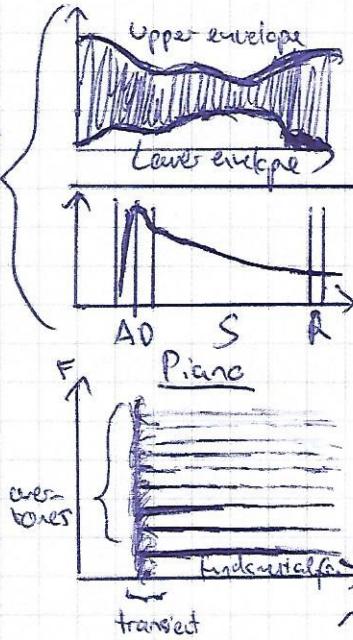


- pitch perception is logarithmic: $F(p) = 2^{\frac{p-69}{12}} \cdot 440\text{Hz}$
- loudness perception is logarithmic as well: $\text{dB}(I) = 10 \cdot \log_{10} \frac{I}{I_{\text{threshold}}}$
 $I_{\text{threshold}} = 10^{-12} \text{W/m}^2$. ← threshold of hearing
 $I_{\text{per}} = 10^{-10} \text{W/m}^2$. ← threshold of pain
 $\rightarrow \text{dB}(2I) \approx \text{dB}(I) + 3$ → double intensity $\hat{=} +3\text{ dB}$
 $\rightarrow \underline{\text{phon}} = \text{perceived loudness}$

- dynamics & timbre:

ADSR model

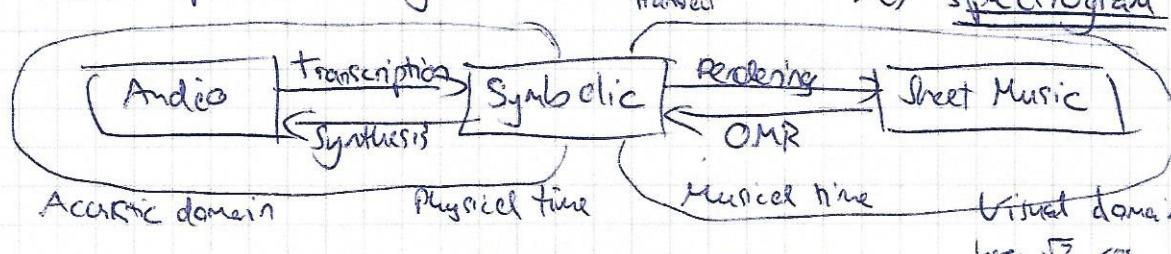
Attack — noise, transients
Decay
Sustain
Release



- variations in amplitude/frequency:

- tremolo
 $\hat{=} \text{amplitude modulation}$
- vibrato
 $\hat{=} \text{frequency modulation}$

- OMR = Optical Music Recognition

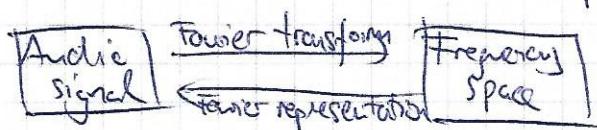


Chapter 2. Fourier Analysis of Signals

idea: - compare signal with sinusoids of various frequencies $w \in \mathbb{R}$
 - for each $w \in \mathbb{R}$, we get a magnitude coefficient $a_w \in \mathbb{R}_{\geq 0}$ and a phase coefficient $\varphi_w \in \mathbb{R}$

$$\text{sinusoid} := A \sin(2\pi(wt - \varphi))$$

↑
amplitude frequency phase



How to compose two signals f.g.:
 $\int f(t) \cdot g(t) dt$

"Fourier transform": $a_w := \max_{t \in [0,1]} \left| \Phi_{w,\varphi} \right|$
 $\varphi_w := \arg \max_{\varphi \in [0,1]} \left| \Phi_{w,\varphi} \right|$

$$\left. \Phi_{w,\varphi} \right| \text{ where } \Phi_{w,\varphi} := \int_{t \in \mathbb{R}} f(t) \sqrt{2} \cos(2\pi(wt - \varphi)) dt$$

comparing f with a sinusoid

This can be simplified by introducing complex numbers:

$$c_\omega := \underbrace{\frac{d\omega}{JZ}}_{\substack{\text{magnitude} \\ \hat{=} \text{norm of } c_\omega}} \cdot e^{2\pi i (-\rho_\omega)} \in \mathbb{C}$$

phase
 $\hat{=} \text{angle } \gamma$

- Fourier transform of f : $\hat{f}(\omega) := c_\omega$

A main result in Fourier analysis is that

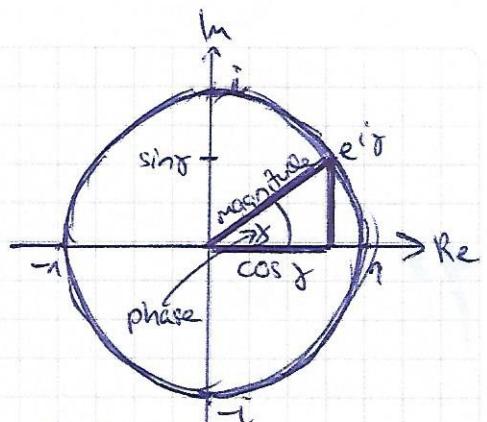
$$\boxed{\hat{f}(\omega) = \int_{t \in \mathbb{R}} f(t) \cdot e^{-2\pi i \omega t} dt}$$

\rightsquigarrow continuous-time (CT) signal

Discrete Fourier Transform (DFT)

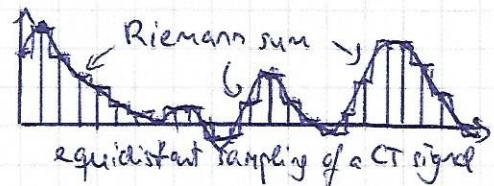
- $F_{S'} := \frac{1}{T}$ (Hz)
Sampling rate Sampling period

Approximation



$$e^{i\gamma} = \cos \gamma + i \sin \gamma$$

(a weighted superposition theory yields
the Fourier representation)



- Nyquist theorem: To reconstruct frequency Ω , the sampling rate must be at least $F_{S'} \geq 2 \cdot \Omega$.
- let $\Omega := F_{S'}/2$, then Ω is the Nyquist frequency and the signal is Ω -limited
- often, $F_{S'} = 44,1$ kHz as $\Omega = 22$ kHz > audible frequency range
- to calculate the Fourier transform, we take a finite amount of samples N and a finite amount of frequency bins M ,
- set $N=M$ to use the Fast Fourier Transform algorithm.

$$\rightsquigarrow \boxed{X(k) = \sum_{n=0}^{N-1} x(n) \cdot e^{-2\pi i kn/N}}$$

$$\text{and } F(k) := \frac{k \cdot F_{S'}}{N \cdot T} = \frac{k \cdot F_{S'}}{N}$$

\rightsquigarrow discrete-time (DT) signal \rightsquigarrow there are as many bins as samples (N)
only consider the lower bins! $\Rightarrow F(N/2) = F_{S'}/2 = \Omega$

Example: $F_{S'} = 44,1$ kHz, then we have 22050 relevant bins (each 1 Hz wide)

Suppose N is a power of two.

In that case, the FFT algorithm can compute the DFT in $O(N \log_2 N)$ instead of $O(N^2)$.

\Rightarrow Always sample with powers of two (window size)?

Short-Time Fourier Transform

Problem: DFT hides time information completely
Idea: consider small "windows"

Let $N \in \mathbb{N}$ be the window size,
 $H \in \mathbb{N}$ the hop size and
 $w: [0, N-1] \rightarrow \mathbb{R}$ a window function.

~ The discrete STFT of the signal x with the k^{th} Fourier coefficient at k^{th} time slice:

$$\mathcal{Z}(m/k) := \sum_{n=0}^{N-1} x(n+mH) \cdot w(n) \cdot e^{-\frac{2\pi i k n}{N}}$$

with $m \in \mathbb{Z}$, $k \in [0, N/2]$.

The meaning of m and k is as follows:

$$T(m) := \frac{m \cdot H}{FS'} \quad F(k) := \frac{k \cdot FS'}{N}$$

Let F_S be the sample rate and N the window size.

Then there are $\frac{N}{2}$ frequency bins, each F_s/N Hz apart.
 (From one equally spaced upto the Nyquist Frequency.)

Tradeoff

deoff: good time resolution \checkmark^N = bad frequency $\text{FB}(N)$

- width size N proportional to the number of bins $N/2$ \hookrightarrow bins for $\Omega/2!$
 - width size N anti-proportional to width of bins $F_D^x / N = \frac{\Delta}{\text{bins}} = \frac{fs/2}{N/2}$

Example: $N = 4096$ $H = 2048$
 $f_S = 44.1 \text{ kHz}$

Suppose there are 40 flowers.

Since there are 40 fewer $\frac{1}{2}$ -pounds than 1-pounds, there are 20 less.

~~104E~~ 1 PC 11-16-2

Then the signal is -1140

Then the total is $-1(40) = -10 \cdot 44.1 \approx 1,16$ seconds long.

Chapter 3. Music Synchronisation

↳ a procedure which, for a given position in one representation of a piece of music, determines the corresponding position in another representation (or interpretation).

Two steps: (mid-level)

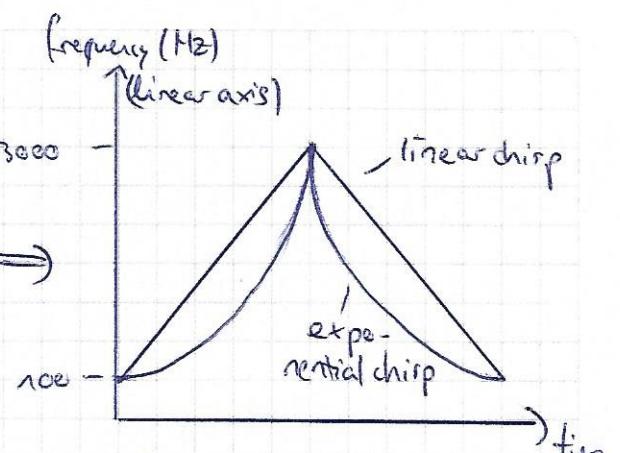
- Feature representation (e.g. Chroma, MFCC, ...)
 - Dynamic time warping (DTW) to align sequences



Exercise 1/2

- chirp is a continuously increasing in frequency.
- linear: the function that describes the frequency is linear: $(f_0 + \frac{1}{2} \cdot t) \cdot t$ \Rightarrow
- exponential: the derivative is exponential

The exponential chirp starts off slower but gets much faster towards the end. Because of the logarithmic perception of pitch, an exponential chirp may be perceived as a more consistent/even increase in pitch.



- There are several ways to combine signals:
Let s_1, s_2 be pure sinusoids with frequency f_1, f_2 respectively. The signal
 - $s_1 + s_2$ is the superposition of both pitches (polyphony) (contains f_1, f_2)
 - $s_1 \cdot s_2$ performs ring modulation ($f_1 \pm f_2, f_1 + f_2$), which changes the timbre
 - $\sin(f_1 \cdot 2\pi t + \frac{1}{f_2} s_2(t))$ performs a frequency modulation of f_2 into f_1
- More complex combinations are possible, like a sawtooth wave: (fundamental frequency f_0)

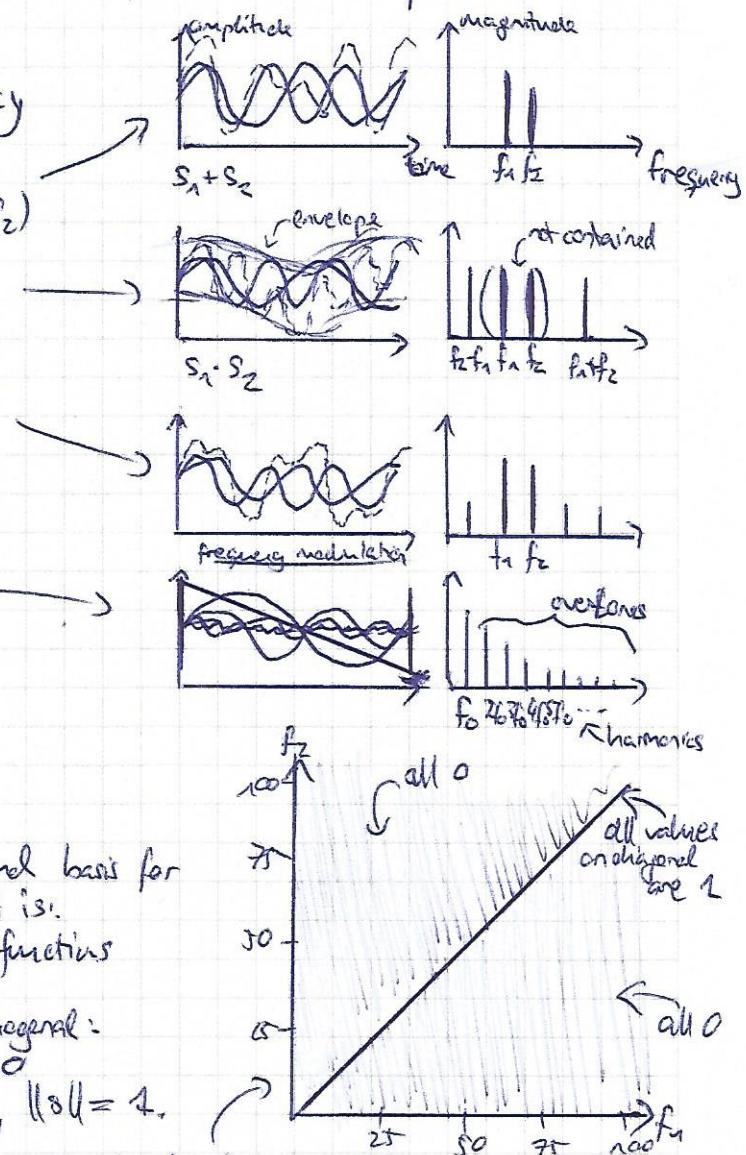
$$s(x) = \sum_{n=1}^{\infty} \frac{1}{n} \cdot \sin(2\pi f_0 x)$$

- Why can every real-valued function be decomposed into Taylor space?

Because Taylor space has all sinusoids as an orthonormal basis. In other words,

$B = \{ \sin(2\pi w t) \mid w \in \mathbb{R} \}$ is an orthonormal basis for the ∞ -dimensional Fourier space, that is:

- B spans the space of real-valued functions
- B is minimal in doing so
- two sinusoids s_1, s_2 are orthogonal: $\langle s_1, s_2 \rangle = \int s_1(t) \cdot s_2(t) dt = 0$
- all sinusoids are normalized, i.e., $\|s\| = 1$.



This can be visualized by plotting the inner product of different sinusoids.

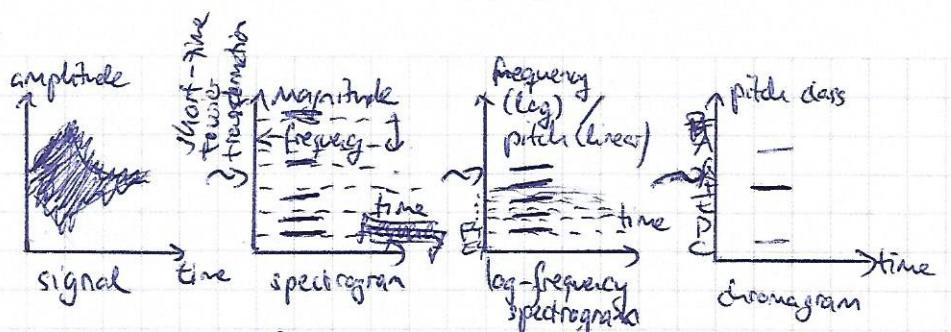
\Rightarrow The Fourier transform is just a change of basis in the space of real-valued functions.



3.1 Chromograms

- calculate the STFT
- do a logarithmic compression to have an axis linear in pitch, this uses the bandwidth of a frequency bin:

(cutoff frequencies) $BW(p) = F(p+0.5) - F(p-0.5)$
(pitch ± 50 cent)



3.2 Dynamic Time Warping

Idea: Compare two time-dependent sequences to find a good local alignment.

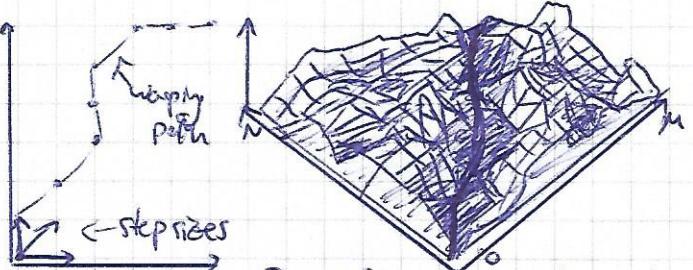
Given two sequences of feature vectors, calculate a cost matrix C that contains the cosine similarities between all feature vectors.

A warping path is a path from the bottom-left to the top-right with specific allowed step sizes.

- noncausal: no going backwards
 - boundary condition: start/end point
- It accumulates the cost of all passed cells.

We aim to minimize the cost.

DTW distance := cost of optimal warping path



3D surface of the cost matrix with dynamic warping path

As the number of warping paths is exponential, we use dynamic programming instead ($O(NM)$). We calculate the accumulated cost matrix D as follows:

$D(n, 1) = \sum_{k=1}^m C(k, 1)$ $D(1, m) = \sum_{k=1}^n C(1, k)$ $D(n, m) = C(n, m) + \min \left\{ \begin{array}{l} D(n-1, m-1) \\ D(n-1, m) \\ D(n, m-1) \end{array} \right\}$

=)
depends
on step sizes

$D(n, 1)$	
$D(n-1, m-1)$	$D(n, m)$
$D(n-1, m)$	$D(n, m)$
$D(n, m-1)$	$D(n, m)$

↳ does not guarantee optimal path!

Each cell is the DTW distance for a subsequence of the given sequences.

To determine the optimal path, backtrack from the upper right corner by choosing the minimum of all possible origin cells.

* alternatively, set $D(n, 0) = D(0, n) = +\infty$

Step sizes can be changed to account for too large temporal deformations (e.g., \uparrow):

\nwarrow	C	D
\nearrow	$\begin{matrix} 1 & 1 & 1 & 7 & 0 \\ 6 & 8 & 8 & 0 & 1 \\ 1 & 3 & 3 & 5 & 4 \\ 1 & 3 & 3 & 5 & 4 \\ 1 & 1 & 1 & 7 & 6 \end{matrix}$	$\begin{matrix} 10 & 10 & 11 & 17 & 10 \\ 9 & 11 & 12 & 17 & 8 \\ 3 & 5 & 12 & 12 & 12 \\ 2 & 4 & 5 & 8 & 12 \\ 1 & 2 & 3 & 10 & 10 \end{matrix}$
\searrow		Fabs. diff. $(2, 0, 0, 8, 2)$
\uparrow		optimal path



- Variations of DTW:
- step sizes to avoid temporal deformation
 - local weights to favor diagonal movement
 - global constraints to filter "safe" paths =>

Note that an optimal path does not necessarily exist (esp. for sequences of different length).

3.3 Applications

→ robustness vs. expressiveness of features

Chromagen: → robust to variations in timbre and dynamics

- characterize melodic, harmonic and possibly temporal progression

multimedial music navigation: Interpretation, Score Viewer

performance analysis with tempo curves: find commonalities & characteristics of interpretations

→ calculate DTW with reference version (annotated with BPM)

online score following; automatic accompaniment

Chapter 4 Music Structure Analysis

4.1 General Principles

Goal: Divide a given music representation into temporal segments that correspond to musical parts.
(motifs, phrases, sections, ...)

→ methods are usually repetition-, novelty- or heterogeneity-based.

→ Waveform - loudness
Recurrence - structure

Different features are used in different contexts: e.g.,

Also, long window sizes in STFT can smooth out instant variations.
→ homogeneity-based segmentation



A B C B

↓

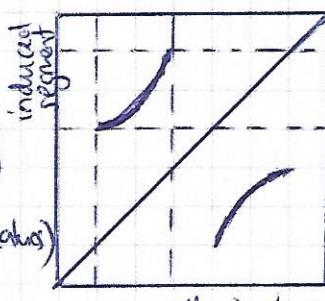
- Chroma - harmonics, but no timbre
- MFCC - timbre (envelope shape)
- Tempo - rhythm, but no other features

4.2 Self-Similarity Matrices

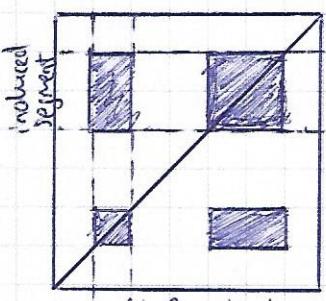
Ideas: pairwise comparison of all feature vectors, yielding a SSM (self-sim. matrix).
(Compared to example with cosine similarity.)

Structures that appear are:

- blocks ≈ homogeneity (consistent features)
- paths ≈ repetition of subsequences



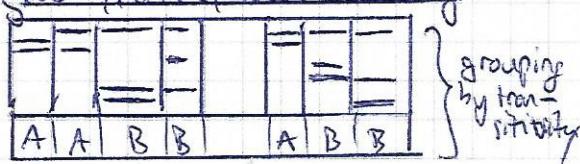
a path structure



a block structure

Enhancement strategies:

- masking, downsampling
- path smoothing: averaging / low-pass filter in the direction of the main diagonal
- tempo-invariant smoothing, forward/backward smoothing → (assuming no tempo differences)
- transition-invariant self-similarity matrix
- global/relative/local thickening



$$S_{\text{LT}}(n, m) = \max_{0 \leq \theta \leq 1} \left[\frac{1}{L} \sum_{l=0}^{L-1} S(n+l, m+l \cdot \theta) \right]$$

$$S^{\text{TI}}(n, m) = \max_{t \in [0, 1]} [S(P^t(n), m)]$$

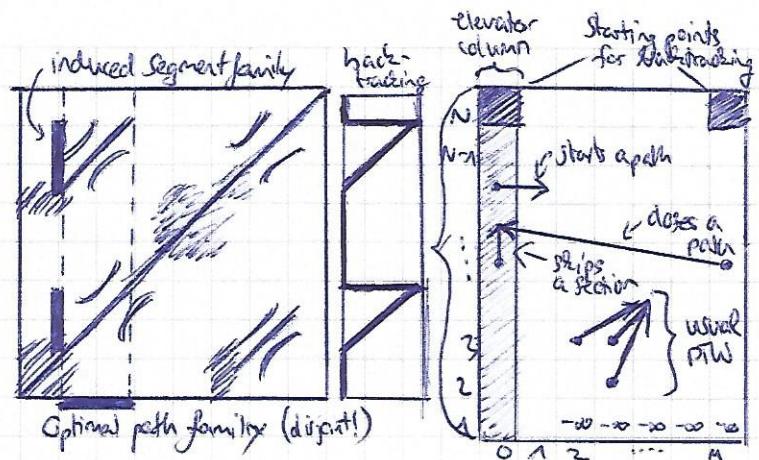


4.3 Audio thumbnailing

Goal: obtain the most representative section of a recording (e.g., clowns).

Idea: introduce fitness measure for a given segment that describes how well and how much a segment describes a recording; maximize the fitness.

induced segment family:
non-overlapping repetitions of a given segment



Try to align the segment with several induced segments with a modified DTW procedure.
 ↳ differences: align the segment or multiple times, but completely, potentially not covering the whole recording and maximizing the path family score.
 i.e. skip sections of α and go back to the beginning of α ('carriage return')

The accumulated cost matrix is defined by

$$\begin{aligned} D(n, m) &= S^\alpha(n, m) + \max \{ D(i, j) | (i, j) \in \mathbb{P}(n, m) \} && - \text{step size condition as usual in DTW} \\ D(n, 0) &= \max \{ D(n-1, 0), D(n-1, m) \} && - \text{elevator column: skip segments, close up paths} \\ D(n, 1) &= D(n, 0) + S^\alpha(n, 1) && - \text{start a new path compared coming from the elevator column} \\ D(1, m) &= -\infty \quad (m \geq 2) && - \text{forces the first path to come from the elevator column} \end{aligned}$$

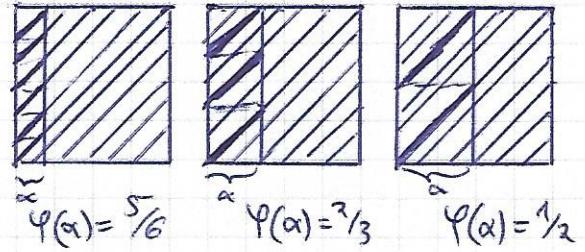
Eliminating the self-explanation, the normalized score describes how well this segment explains other segments.
 The coverage describes how much of the recording it explains; this is the sum of all induced segments.

- short segments explain usually well, but not much
- long segments explain usually much, but not well

We unify this with a harmonic mean (F-measure):

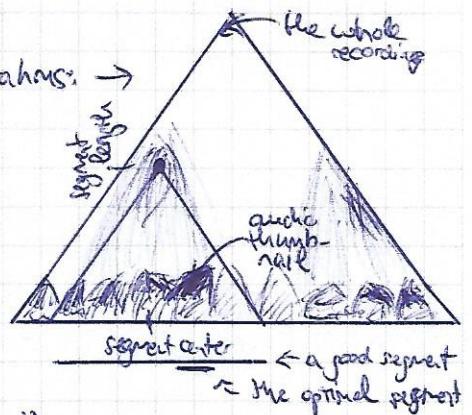
$$F(\alpha) := 2 \cdot \frac{\bar{\alpha}(\alpha) \cdot \bar{\beta}(\alpha)}{\bar{\alpha}(\alpha) + \bar{\beta}(\alpha)} \quad \bar{\alpha}: \text{precision / score} \quad \bar{\beta}: \text{recall / coverage}$$

(favors shorter segments due to self-explanation)



Represent each segment by its length and center.

We can indicate the fitness of all segments in a scope plot. Brahms: →

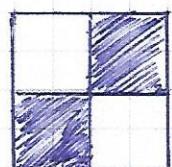


4.4 Novelty-based Segmentation

Goal: locate musical changes in time (segment boundaries).

Novelty Detection: slide a (Gaussian) checkerboard kernel across the diagonal of an SSMatrix and calculate similarity. This detects block-like structures.

→ novelty function (changes in timbre, pitch...)



kernel (size matters!)

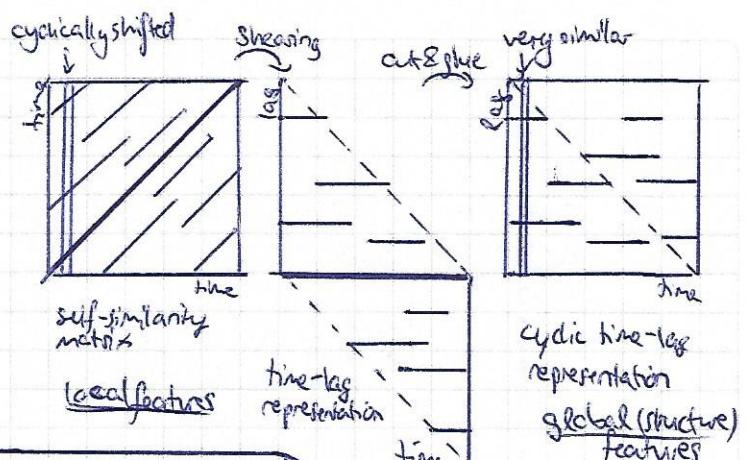


Structure Features (lag-level novelty)
combines local and global characteristics.
frame-wise entire frame

→ time-lag matrices

This transforms frame vectors so that structurally similar frames are no longer cyclically shifted, but almost identical.

→ new (high-level) feature representation
(structure features)

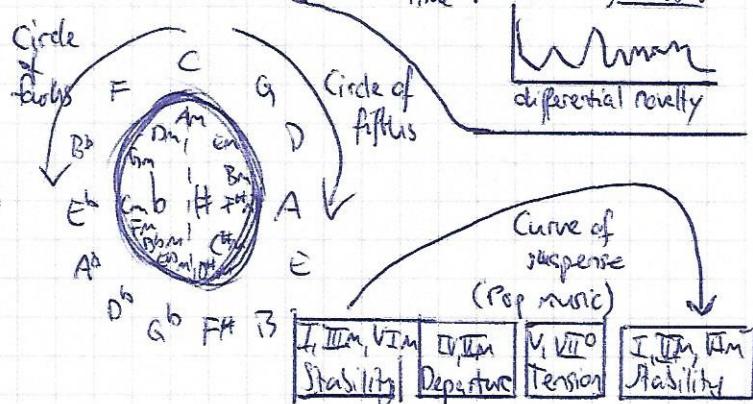


Chapter 5: Chord Recognition

5.1 Basic Theory of Harmony

Intervals, frequency ratios, chords, scales, circle of fifths / fours

functional harmony theory models chord progressions and their fitness.

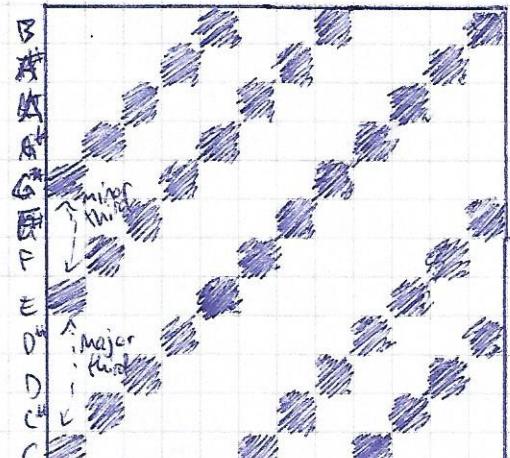


5.2 Template-based Chord Recognition

Based on chroma features, compare frames with given chord triplets (cosine similarity) and find the best match. (Ignores functional harmony theory.)

$$\lambda_n := \underset{\substack{\text{word of n-th frame} \\ \sim \text{template frame}}}{\operatorname{argmax}}_{t \in T} S(t, x_n)$$

- However, there are multiple problems: Some chords are very similar, no functional harmonies are regarded, there are many chord variations (e.g., inversions) (chord ambiguity)
e.g.: $S(t_C, t_C) = S(t_G, t_E)$
- overtones & logarithmic pitches & major-minor-confusion
- tuning issues, broken chords, ...



Templates for Major Triads

Evaluation:

$$F = \frac{Z \cdot P \cdot R}{P + R}$$

$$R = \frac{TP}{TP + FN}$$

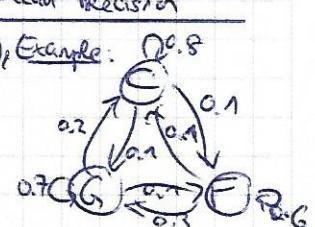


5.3 HMM-Based Chord Recognition

Takes into account the probability of harmonic progression using Hidden Markov Models.

We assume that $P(s_{n+1} = \alpha_j | s_n = \alpha_i, \dots) = P(\alpha_{n+1} = \alpha_j | \alpha_n = \alpha_i, \dots)$. → "Memoryless" (Markov property), Example: 0.8

$$\begin{cases} a_{ij} := P(s_{n+1} = \alpha_j | s_n = \alpha_i), \sum_{j=1}^J a_{ij} = 1 \\ c_i := P(s_1 = \alpha_i) \end{cases}$$





Hidden Markov Models

Goal: uncover relation between observed feature vectors and underlying chords. } Emission probabilities tell whether a chord generates a feature vector.

A discrete HMM is a tuple $\Theta := (\Omega, A, C, \Omega^*, B)$ with

- Ω : set of states α_i
- A : state transition probabilities a_{ij}
- C : initial state probabilities c_i
- Ω^* : set of observation symbols β_k
- B : emission probabilities b_{ik} (that state α_i emits β_k)

A, C , and B are determined by experts or with learning procedures, C and Ω^* are fixed.

Problem: Uncovering the hidden state sequence

Given an observation sequence, we want to assign chord labels that best explain the observations.

Evaluation Problem: How likely is a given observation sequence for a given HMM? (How well does it fit?)

$$P(O, S | \Theta) = c_{i_1} \cdot b_{i_1 k_1} \cdot a_{i_1 i_2} \cdot b_{i_2 k_2} \cdots a_{i_n i_1} \cdot c_{i_1} \cdot b_{i_1 k_1}$$

$$P(O | \Theta) = \sum_{S=(s_1, \dots, s_N)} P(O, S | \Theta) \quad - \text{ sum over all possible explaining state sequences (one of these is optimal!)}$$

Uncovering Problem: Given Θ and O , what is the single state sequence that "best explains" O ? $S = (s_1, \dots, s_N)$

$$\begin{aligned} \text{Prob}^* &= \max_{S=(s_1, \dots, s_N)} P(O, S | \Theta) \\ S^* &= \underset{S=(s_1, \dots, s_N)}{\operatorname{argmax}} P(O, S | \Theta) \end{aligned} \quad \left\{ \begin{array}{l} \text{the state sequence that contributes most to } P(O | \Theta), \\ \text{the probability Prob}^* \text{ may be very small!} \end{array} \right.$$

Viterbi algorithm: fill matrix similar to DPP

$$P(i, 1) = c_i \cdot b_{i k_1}$$

$$D(i, n) = b_{i k_n} \cdot \max_{j \in \Omega, i \neq j} (a_{ji} \cdot D(j, n-1))$$

$$\text{Prob}^* = \max_{i \in \Omega} D(i, N)$$

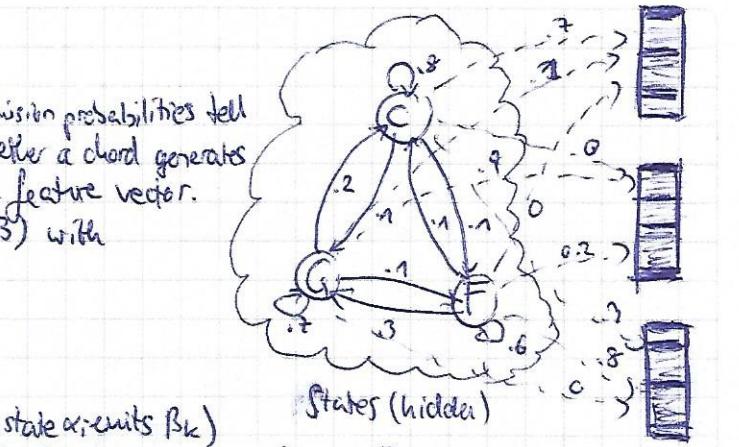
Application to Chord Recognition

$$\Omega = \{C, C^+, \dots, B, C_m, C^{\#}_m, \dots, B_N\}$$

As for ω : Certain chord transitions are

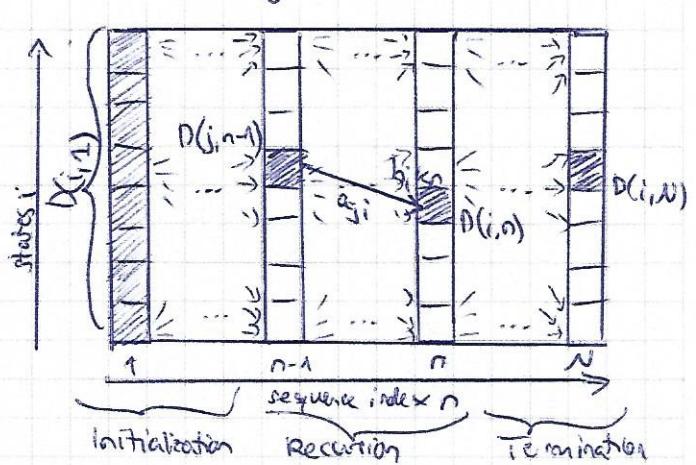
more likely than others (change from a tonic to a dominant is much more likely than transposing by one semitone).

In practice: Training - Baum-Welch algorithm,
Evaluation - Viterbi algorithm

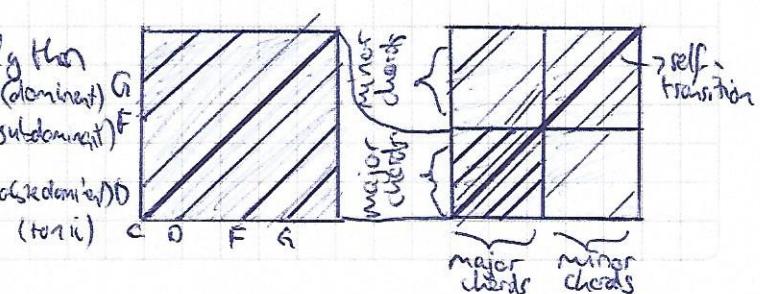


Observation
(visible)

States (hidden)



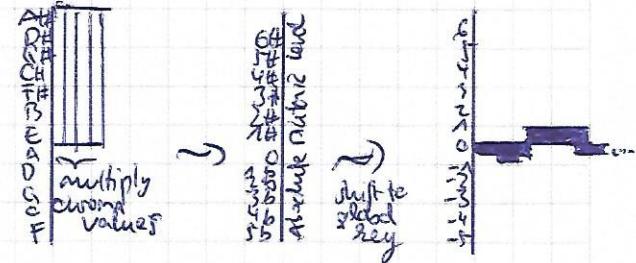
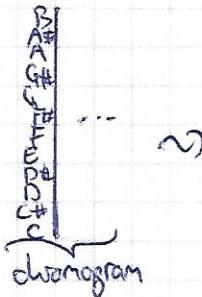
Initialization Recursion Termination





Local Key Detection

- detect modulations in key
- chroma smoothing
- reorder to perfect fifths (circle of fifths!)
- multiply chroma values
- shift to global key



Chapter 6: Tempo & Beat Tracking

Given an audio recording, determine the periodic sequence of beat positions ("tapping the feet when listening to music").

6.1 Onset Detection

- onset = instant that marks the beginning of a transient
- ↳ percussive music: energy-based novelty
 - ↳ nonpercussive music: spectral-based approach

- Steps:
- feature representation that reflects onsets better
 - apply derivative operator to derive novelty function
 - peak-picking algorithm

} see novelty in music structure analysis

Amplitude envelope

Attack | Decay

Onset Transient → unpredictable/chaotic

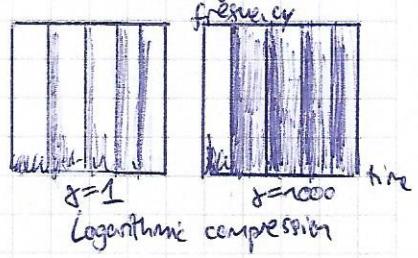
→ Measure
→ Tactus (beat)
→ Tatum (atom)

Energy-based novelty: - shift a bell-shaped function along the signal (local energy)

- take discrete derivative
- half-wave rectification (thresholding)
- (optional) - logarithmic compression to account for logarithmic perception of loudness

Spectral-based novelty:

- logarithmic compression
- derivative, half-wave rectification
- accumulation into novelty function
- normalization (subtract local average)



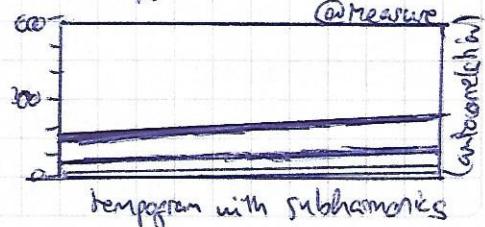
6.2 Tempo Analysis

extraction of tempo information is difficult due to tempo rubato (freely speeds up and down) and syncopation (stressing notes outside the grid of beat positions)

- beat: steady pulse that drives beat forward
- tempo: speed of the pulse

~ strategy: analyze the novelty curve with respect to reoccurrence or quasi-periodic patterns; avoid peak picking!

A tempogram is a time-tempo representation that encodes the local tempo of a music signal over time.





Fourier Tempogram =
 Windauer sin- usodal → - calculate a novelty curve or score (Δ)
 → compute spectrogram (STFT) of Δ
 → convert frequency axis into tempo

Autocorrelation Tempogram:
 ↳ sliding inner product - compare Δ with time-delayed versions of itself
 → convert lag-axis into tempo axis

$\text{Lag} = 0$: $\text{Lag} = 0.2$: $\text{Lag} = 0.4$:

Fourier ("fast") vs. Autocorrelation ("slow")
 - reveals Δ 's periodicities
 - emphasizes harmonics
 - @ tempo / @ tactus
 - reveals Δ 's self-similarities
 - emphasizes subharmonics
 - @ tactus / @ measure

6.3 Beat & Pulse Tracking

Goal: Given the tempo, find best sequence of beats.

PLP is a periodic enhancement of the novelty curve.

The accumulator introduces error subwindows and allows peak picking. The locality of the samples allows for tempo variations.

(Optionally, additional information (rough tempo range) can be used.)
 However, the predominant pulse level can change (maximum!).

→ The detected beat positions can then be used as adaptive windows for improved feature extraction.

Chapter 7. Content-Based Audio Retrieval

Query-by-Example: { Audio identification / fingerprinting
 Audio matching
 Cover song identification

Audio identification / fingerprinting

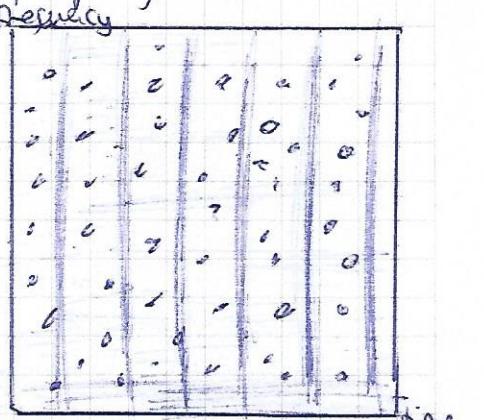
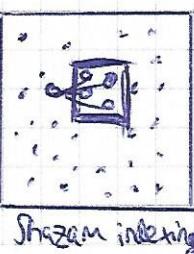
Goal: Given a short audio query fragment, identify the original audio recording the query is taken from. (Shazam)

Idea: An audio fingerprint is a content-based compact signature that summarizes some specific audio contents

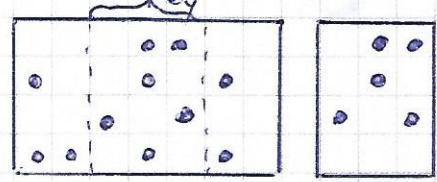
- constellation map of a spectrogram
- slide constellation map for query over the database
- requires clever indexing strategies

Idea: use pairs of peaks to increase hash specificity

- fix anchor point
- define target zone
- use pairs of peaks ($f_1, f_2, \Delta t$)



Constellation map
query



database

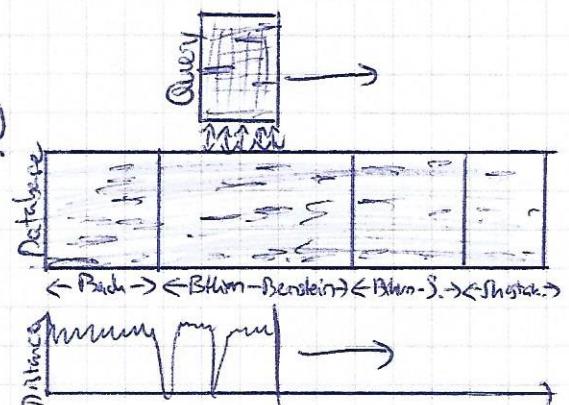
query



Audio matching

Goal: Given a short query audio fragment, find all corresponding audio fragments of similar musical content (piece of music).
 ↗ smoothing, downsampling

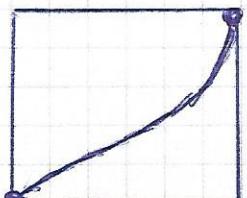
- ↳ use chroma features (allows to find different interpretations, instrumentations, rhythms, ...)
- ↳ then use a matching procedure on the database
- ↳ for tempo differences, use local time warping or multiple scaling
- chroma: inv. to timbre, normalized; inv. to dynamics, smoothing: inv. to local time deviations



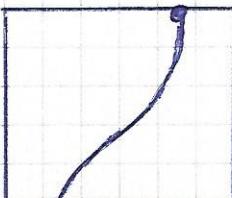
Cover song identification

Goal: Find all kinds of versions, interpretations or cover songs of a piece of music. (potentially radically different interpretations!) ↗ but usually at least one thing doesn't change.

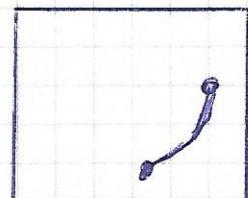
- calculate chroma sequences and a similarity metric
- find a longest common subsequence with retrieval score (local sequence alignment), e.g., Smith-Waterman algorithm used in bioinformatics for nucleotide/protein sequences



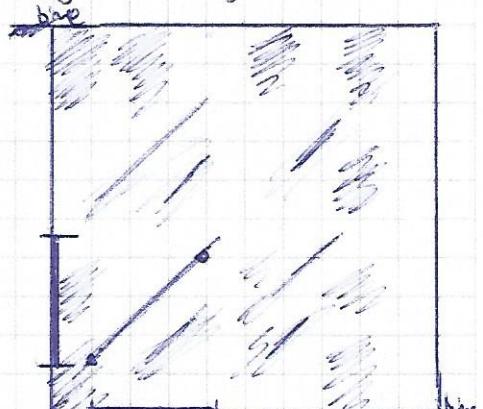
Classical DTW
compare whole recordings
(audio identifiers)



Subsequence DTW
compare query with
database (matching)



Local Alignment
compare most representative
parts (cover song identified!)



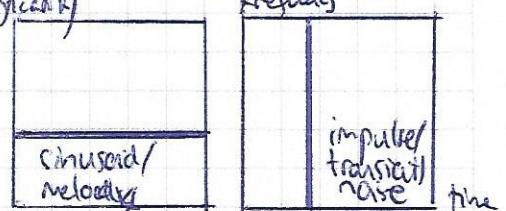
local alignment of a cover song
frequency

Chapter 8. Musically Informed Audio Decomposition

Source separation = decomposition of audio stream into different sources

↗ harmonic-percussive separation

use the fact that percussive elements are vertical lines in a spectrogram and harmonic elements are horizontals



Harmonic: singing voice &

accompaniment

Residual: vibrato/glissando
noise

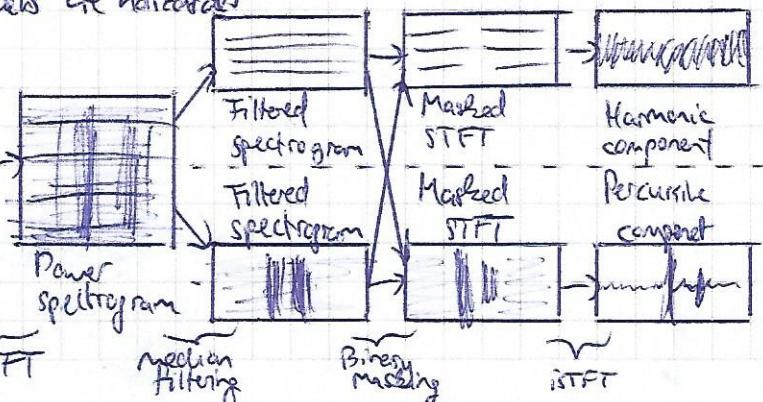
Percussive: drum hits

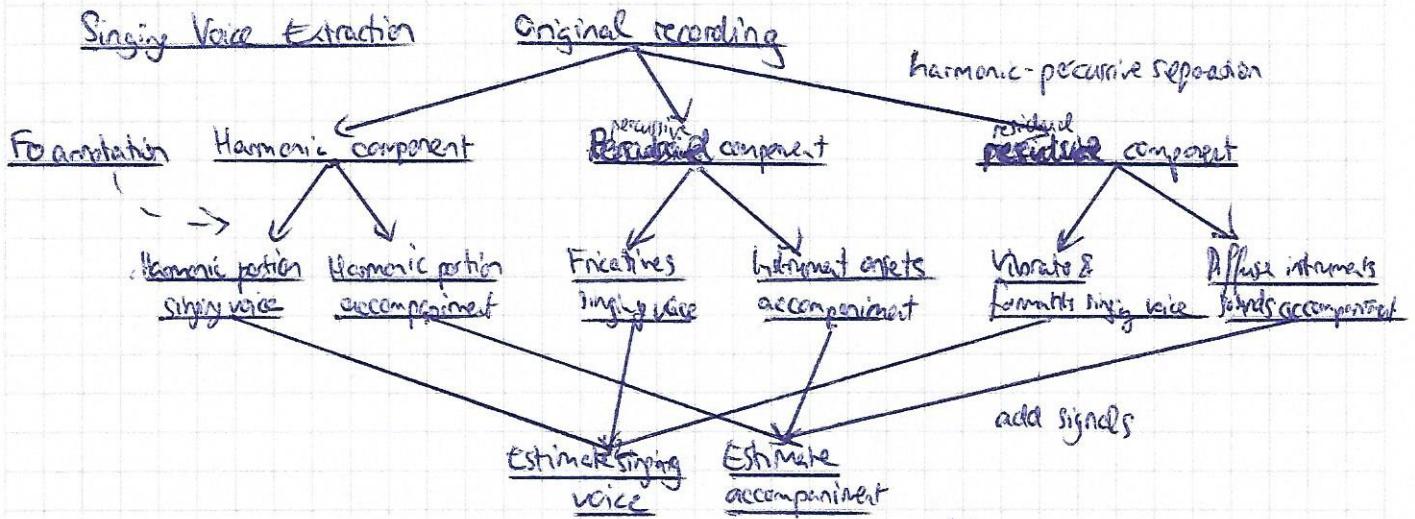
fricatives/hisses

singing voice



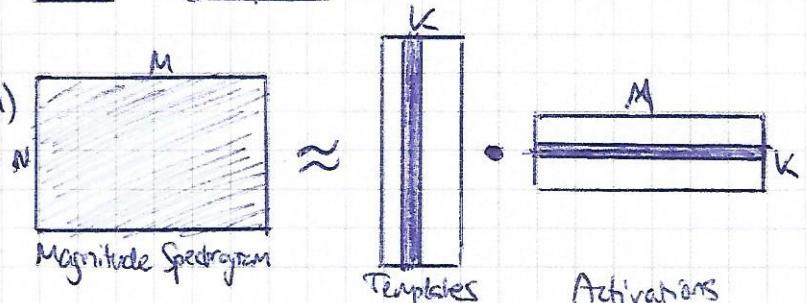
Input
signal





Nonnegative Matrix Factorization

(source-informed audio decomposition)
 factorization of ce nonnegative
 matrix into two nonnegative
 matrices



- templates: pitch + timbre,
 ("how does it sound")

- activations: onset time + duration ("when does it sound")

→ Both matrices are initialized (e.g., randomly) and adapted via a learning procedure.
 ↳ or musically informed: prototypical pitch features

Audio Mosaicing

Target signal: "let it Bee" (Beetles)
 Source signal: "Bees buzzing"
 Mosaic signal: "let it Bee"

} Idea: replace activation matrix with bee activations

Bzzzzz....

LE FIN □