

Audio Features



Fourier Transform

- Tells **which** notes (frequencies) are played, but does not tell **when** the notes are played
- Frequency information is averaged over the entire time interval
- Time information is hidden in the phase

→ Windowed Fourier Transform (WFT)
Short Time Fourier Transform (STFT)
(Dennis Gabor, 1946)

Short Time Fourier Transform

Idea:

- To recover time information, only a **small section** of the signal is used for the spectral analysis
- This section is determined by a **window function** $g : \mathbb{R} \rightarrow \mathbb{R}$ ($g \in L^2(\mathbb{R}), \|g\| = 1$)

Definition:

STFT w.r.t. g of a signal $f : \mathbb{R} \rightarrow \mathbb{R}$

$$\tilde{f}(\omega, t) := \int_{\mathbb{R}} f(u) \bar{g}(u-t) e^{-2\pi i \omega u} du = \langle f | g_{\omega, t} \rangle$$

with $g_{\omega, t}(u) := e^{2\pi i \omega u} g(u-t), u \in \mathbb{R}$

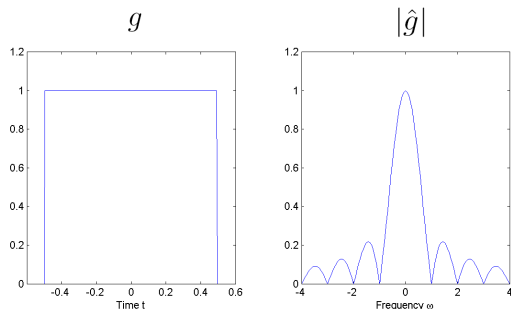
Short Time Fourier Transform

Interpretation:

- $g_{\omega, t}$ represents a „musical note“ of frequency ω which oscillates within the translated window given by $u \rightarrow g(u-t)$
- Inner product $\langle f | g_{\omega, t} \rangle$ measures the correlation between the signal f and the musical note $g_{\omega, t}$

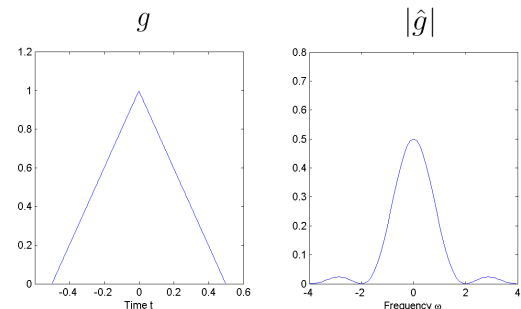
Short Time Fourier Transform

Box window: discontinuities at window boundaries cause artefacts in the frequency domain



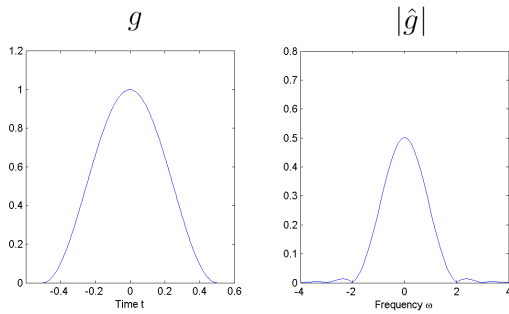
Short Time Fourier Transform

Triangle window



Short Time Fourier Transform

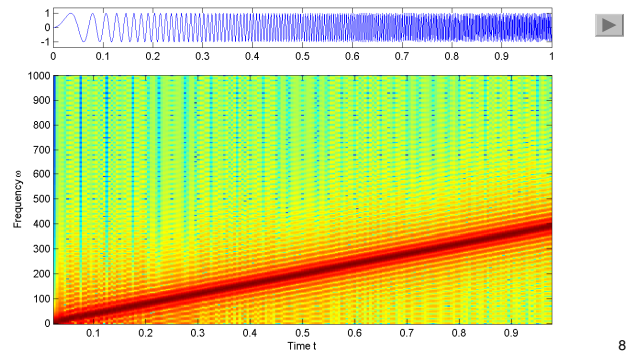
Hann window



7

Short Time Fourier Transform

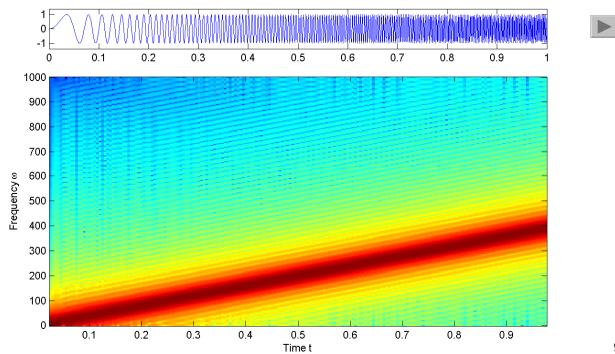
Chirp signal and STFT with **box window** of length 0.05



8

Short Time Fourier Transform

Chirp signal and STFT with **hann window** of length 0.05



9

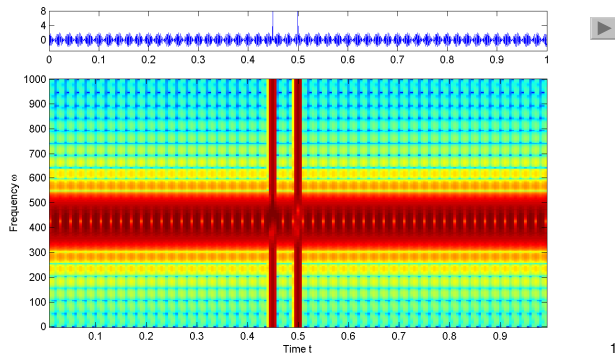
Time-Frequency Localization

- Size of window constitutes a compromise between time resolution and frequency resolution:
 - Large window** : poor time resolution
good frequency resolution
 - Small window** : good time resolution
poor frequency resolution
- Heisenberg Uncertainty Principle: there is no window function that localizes in time and frequency with arbitrary position.

10

Short Time Fourier Transform

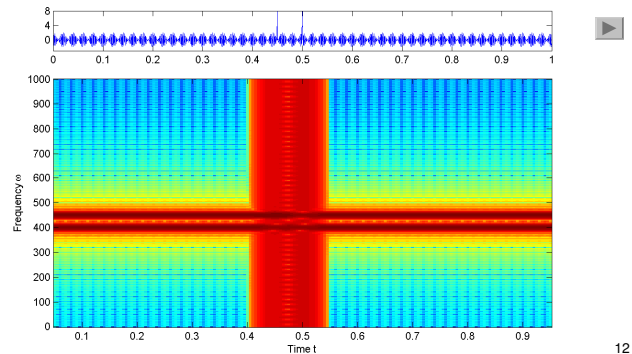
Signal and STFT with hann window of **length 0.02**



11

Short Time Fourier Transform

Signal and STFT with hann window of **length 0.1**



12

Heisenberg Uncertainty Principle

Window function $g \in L^2(\mathbb{R})$ with $\|g\| = 1$

Center

Width

$$t_0 = t_0(g) := \int_{-\infty}^{\infty} t |g(t)|^2 dt \quad T(g) := \left(\int_{-\infty}^{\infty} (t - t_0)^2 |g(t)|^2 dt \right)^{\frac{1}{2}}$$

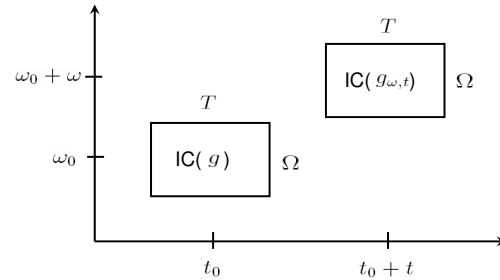
$$\omega_0 = \omega_0(g) := \int_{-\infty}^{\infty} \omega |\hat{g}(\omega)|^2 d\omega \quad \Omega(g) := \left(\int_{-\infty}^{\infty} (\omega - \omega_0)^2 |\hat{g}(\omega)|^2 d\omega \right)^{\frac{1}{2}}$$

$$T(g) \cdot \Omega(g) \geq \frac{1}{4\pi}$$

13

Information Cells

$g_{\omega,t}(u) := e^{2\pi i \omega u} g(u - t)$ "musical note"



14

MATLAB

- MATLAB function SPECTROGRAM
- N = window length (in samples)
- M = overlap (usually $N/2$)
- Compute DFT_N for every windowed section
- Keep lower $N/2$ Fourier coefficients

→ Sequence of spectral vectors
(for each window a vector of dimension $N/2$)

15

Example

Let x be a DT-Signal $x(n) = f(Tn)$

Sampling rate: $1/T = 22050$ Hz

Window length: $N = 4096$

Overlap: $N/2 = 2048$

Hopsize: window length – overlap

Let $v_0 := (x(0), x(1), \dots, x(4095))$

$v_1 := (x(2048), \dots, x(6143))$

$v_2 := (x(4096), \dots, x(8191))$

⋮

v_m corresponds to window $[m \cdot 2048 : m \cdot 2048 + 4095]$

16

Example

Time resolution:

$$\frac{\text{hopsize}}{\text{sampling rate}} = \frac{4096 - 2048}{22050} = 0.093 = 93 \text{ ms}$$

Frequency resolution:

$$v = v_0, \hat{v} := \text{DFT}_N(v)$$

$$\hat{v}(k) \approx \frac{1}{T} \cdot \hat{f} \left(\frac{k}{N} \cdot \frac{1}{T} \right)$$

$$\omega = \frac{k}{N} \cdot \frac{1}{T} = k \cdot \frac{22050}{4096} = k \cdot 5.38 \text{ Hz}$$

17

Pitch Features

Model assumption: Equal – tempered scale

- MIDI pitches: $p \in [1 : 128]$
- Piano notes: $p = 21$ (A0) to $p = 108$ (C8)
- Concert pitch: $p = 69$ (A4) = 440 Hz
- Center frequency: $f_{\text{MIDI}}(p) = 2^{\frac{p-69}{12}} \cdot 440$

→ Logarithmic frequency distribution

Octave: doubling of frequency

18

Pitch Features

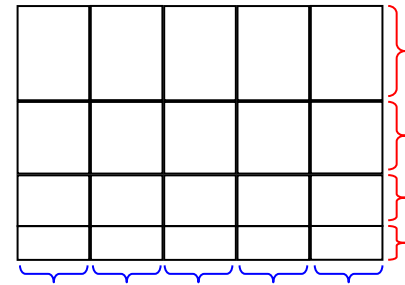
Idea: Binning of Fourier coefficients

Divide up the frequency axis into logarithmically spaced „pitch regions“ and combine **spectral coefficients** of each region to a single **pitch coefficient**.

19

Pitch Features

Time-frequency representation



Windowing in the time domain

Windowing in the frequency domain

20

Pitch Features

Note	MIDI pitch	Center [Hz]	Left [Hz] boundary	Right [Hz] boundary	Width [Hz]
A3	57	220.0	213.7	226.4	12.7
A#3	58	233.1	226.4	239.9	13.5
B3	59	246.9	239.9	254.2	14.3
C4	60	261.6	254.2	269.3	15.1
C#4	61	277.2	269.3	285.3	16.0
D4	62	293.7	285.3	302.3	17.0
D#4	63	311.1	302.3	320.2	18.0
E4	64	329.6	320.2	339.3	19.0
F4	65	349.2	339.3	359.5	20.2
F#4	66	370.0	359.5	380.8	21.4
G4	67	392.0	380.8	403.5	22.6
G#4	68	415.3	403.5	427.5	24.0
A4	69	440.0	427.5	452.9	25.4

21

Pitch Features

Details:

- Let \hat{v} be a spectral vector obtained from a spectrogram w.r.t. a sampling rate $1/T$ and a window length N . The spectral coefficient $\hat{v}(k)$ corresponds to the frequency

$$f_{\text{coeff}}(k) := \frac{k}{N} \cdot \frac{1}{T}$$

- Let $S(p) := \{k : f_{\text{MIDI}}(p - 0.5) \leq f_{\text{coeff}}(k) < f_{\text{MIDI}}(p + 0.5)\}$ be the set of coefficients assigned to a pitch $p \in [1 : 128]$. Then the pitch coefficient $P(p)$ is defined as

$$P(p) := \sum_{k \in S(p)} |\hat{v}(k)|^2$$

22

Pitch Features

Example: A4, $p = 69$

- Center frequency: $f(p = 69) = 2^{\frac{0}{12}} \cdot 440 = 440 \text{ Hz}$
- Lower bound: $f(p = 68.5) = 2^{\frac{-0.5}{12}} \cdot 440 = 427.5 \text{ Hz}$
- Upper bound: $f(p = 69.5) = 2^{\frac{0.5}{12}} \cdot 440 = 452.9 \text{ Hz}$
- STFT with $N = 4096, 1/T = 22050$

\vdots
 $f(k = 79) = 425.3 \text{ Hz}$
 $f(k = 80) = 430.7 \text{ Hz}$
 $f(k = 81) = 436.0 \text{ Hz}$
 $f(k = 82) = 441.4 \text{ Hz}$
 $f(k = 83) = 446.8 \text{ Hz}$
 $f(k = 84) = 452.2 \text{ Hz}$
 $f(k = 85) = 457.6 \text{ Hz}$
 \vdots

23

Pitch Features

Example: A4, $p = 69$

- Center frequency: $f(p = 69) = 2^{\frac{0}{12}} \cdot 440 = 440 \text{ Hz}$
- Lower bound: $f(p = 68.5) = 2^{\frac{-0.5}{12}} \cdot 440 = 427.5 \text{ Hz}$
- Upper bound: $f(p = 69.5) = 2^{\frac{0.5}{12}} \cdot 440 = 452.9 \text{ Hz}$
- STFT with $N = 4096, 1/T = 22050$

\vdots
 $f(k = 79) = 425.3 \text{ Hz}$
 $f(k = 80) = 430.7 \text{ Hz}$
 $f(k = 81) = 436.0 \text{ Hz}$
 $f(k = 82) = 441.4 \text{ Hz}$
 $f(k = 83) = 446.8 \text{ Hz}$
 $f(k = 84) = 452.2 \text{ Hz}$
 $f(k = 85) = 457.6 \text{ Hz}$
 \vdots

$\left. \begin{array}{l} \vdots \\ f(k = 79) = 425.3 \text{ Hz} \\ f(k = 80) = 430.7 \text{ Hz} \\ f(k = 81) = 436.0 \text{ Hz} \\ f(k = 82) = 441.4 \text{ Hz} \\ f(k = 83) = 446.8 \text{ Hz} \\ f(k = 84) = 452.2 \text{ Hz} \\ f(k = 85) = 457.6 \text{ Hz} \\ \vdots \end{array} \right\} S(p = 69)$

$P(p = 69) = \sum_{k=80}^{84} |\hat{v}(k)|^2$

24

Pitch Features

Note:

- $P \in \mathbb{R}^{128}$
- For some pitches, $S(p)$ may be empty. This particularly holds for low notes corresponding to narrow frequency bands.


→ Linear frequency sampling is problematic!

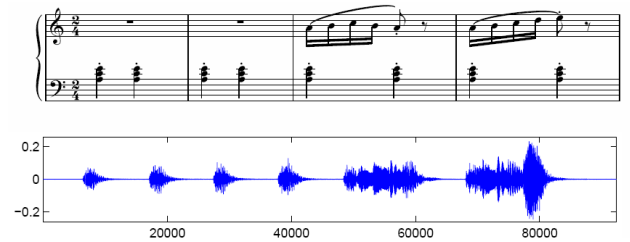
Solution:

Multi-resolution spectrograms or multirate filterbanks

25


Audio Representation

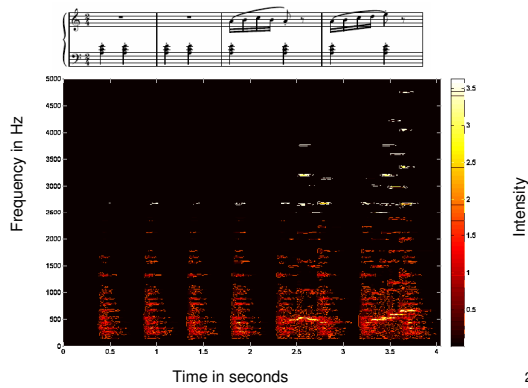
Example: Op. 100, No. 2 by Friedrich Burgmüller 



26


Short Time Fourier Transform

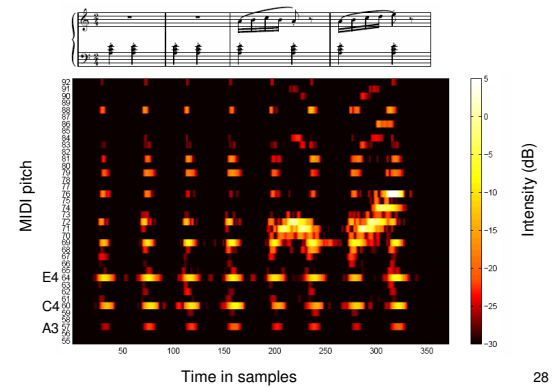
Example: Op. 100, No. 2 by Friedrich Burgmüller 



27


Pitch Features

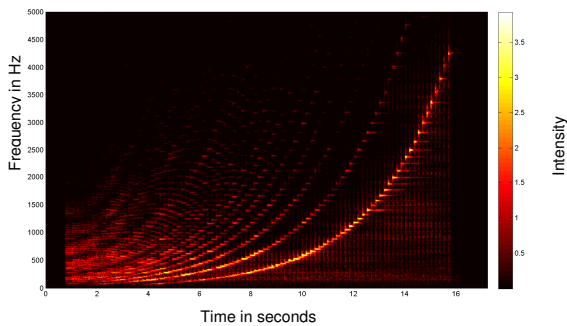
Example: Op. 100, No. 2 by Friedrich Burgmüller 



28


Short Time Fourier Transform

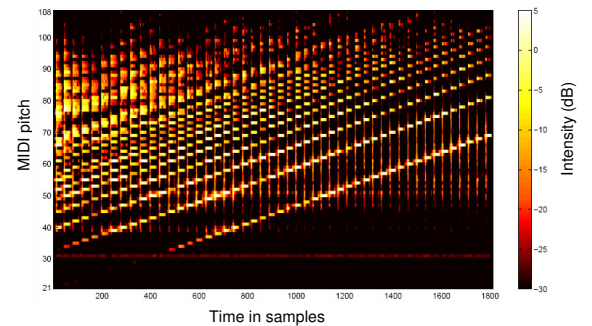
Example: Chromatic Scale 



29

Pitch Features

Example: Chromatic Scale 



30

Chroma Features

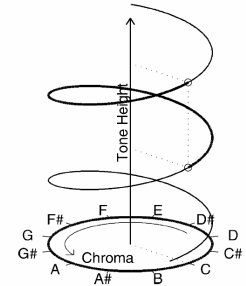
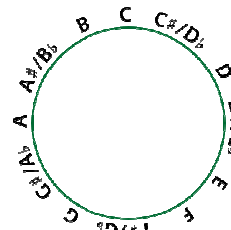
- Human perception of pitch is periodic in the sense that two pitches are perceived as similar in „color“ if they differ by an octave.
- Separate pitch into two components: tone height (octave number) and chroma.
- Chroma : 12 traditional pitch classes of the equal-tempered scale. For example
Chroma C $\cong \{ \dots, C_0, C_1, C_2, C_3, \dots \}$
- Computation: pitch features \rightarrow chroma features
Add up all pitches belonging to the same class
- Result: 12-dimensional chroma vector.

31

Chroma Features

Chromatic circle

Shepard's helix of pitch perception



http://en.wikipedia.org/wiki/Pitch_class_space

Bartsch/Wakefield, IEEE Trans. Multimedia, 2005

32

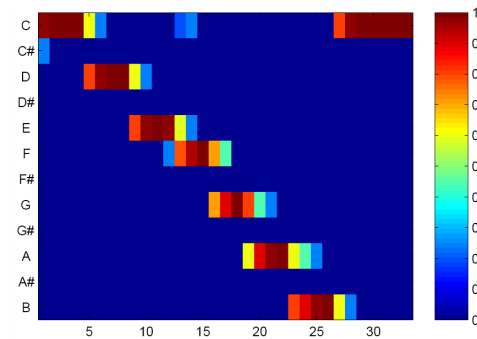
Chroma Features

- Sequence of chroma vectors correlates to the harmonic progression
- Normalization $v \rightarrow \frac{v}{\|v\|}$ makes features invariant to changes in dynamics
- Further quantization and smoothing: CENS features
- Taking logarithm before adding up pitch coefficients accounts for logarithmic sensation of intensity

33

Chroma Features

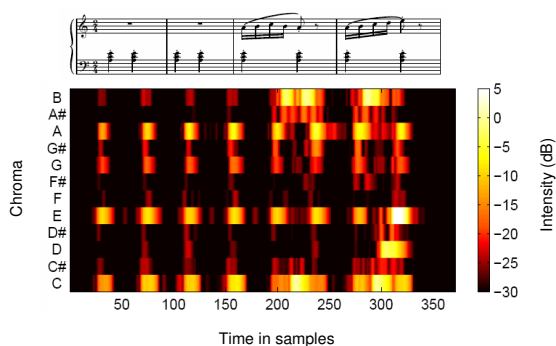
Example: C-Major Scale



34

Chroma Features

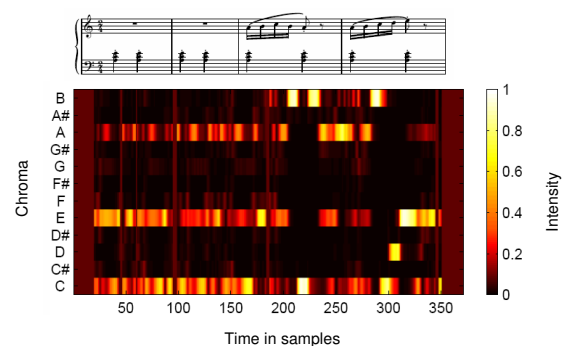
Example: Burgmüller Op. 100, No. 2



35

Chroma Features

Normalization



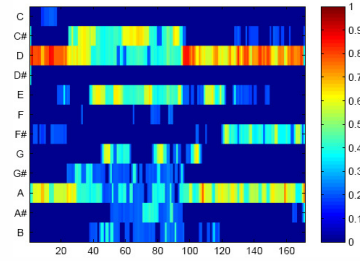
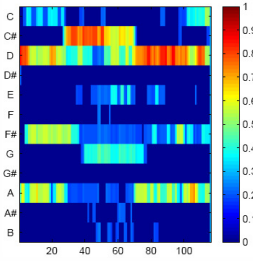
36

Chroma Features

Example: Bach Toccata

Koopman

Ruebsam



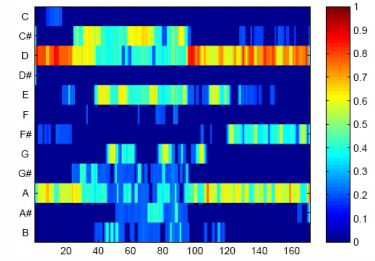
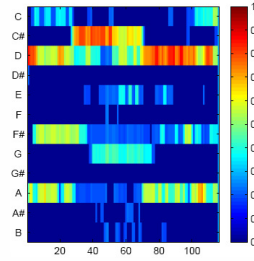
37

Chroma Features

Example: Bach Toccata

Koopman

Ruebsam



Feature resolution: 10 Hz

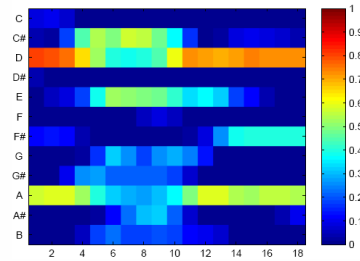
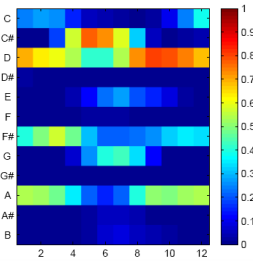
38

Chroma Features

Example: Bach Toccata

Koopman

Ruebsam



Feature resolution: 1 Hz

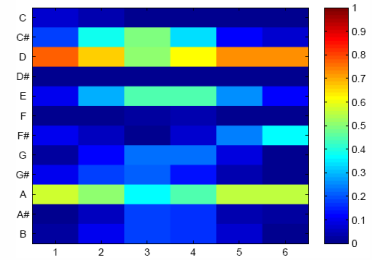
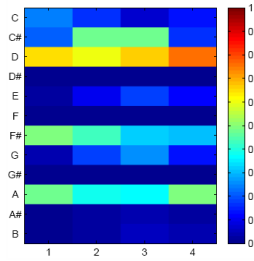
39

Chroma Features

Example: Bach Toccata

Koopman

Ruebsam



Feature resolution: 0.33 Hz

40

Chroma Features

WAV Chroma (10 Hz) GENS (1 Hz)

???

???

???

41

Chroma Features

WAV Chroma (10 Hz) GENS (1 Hz)

Beethoven's Fifth (Bernstein)

???

???

42

Chroma Features

	WAV	Chroma (10 Hz)	CENS (1 Hz)
Beethoven's Fifth (Bernstein)			
Beethoven's Fifth (Piano/Sherbakov)			
???			

43

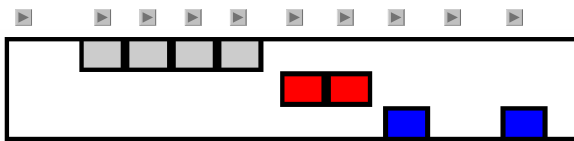
Chroma Features

	WAV	Chroma (10 Hz)	CENS (1 Hz)
Beethoven's Fifth (Bernstein)			
Beethoven's Fifth (Piano/Sherbakov)			
Brahms Hungarian Dance No. 5			

44

Chroma Features

Example: Zager & Evans "In The Year 2525"

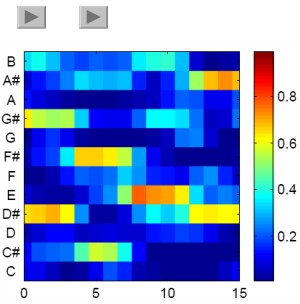


How to deal with transpositions?

45

Chroma Features

Example: Zager & Evans "In The Year 2525"

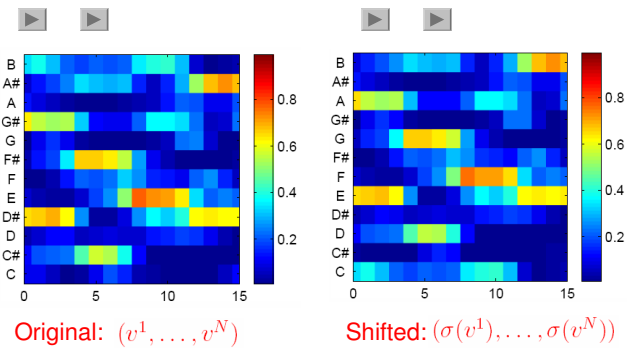


Original: (v^1, \dots, v^N)

46

Chroma Features

Example: Zager & Evans "In The Year 2525"



Original: (v^1, \dots, v^N)

Shifted: $(\sigma(v^1), \dots, \sigma(v^N))$

47