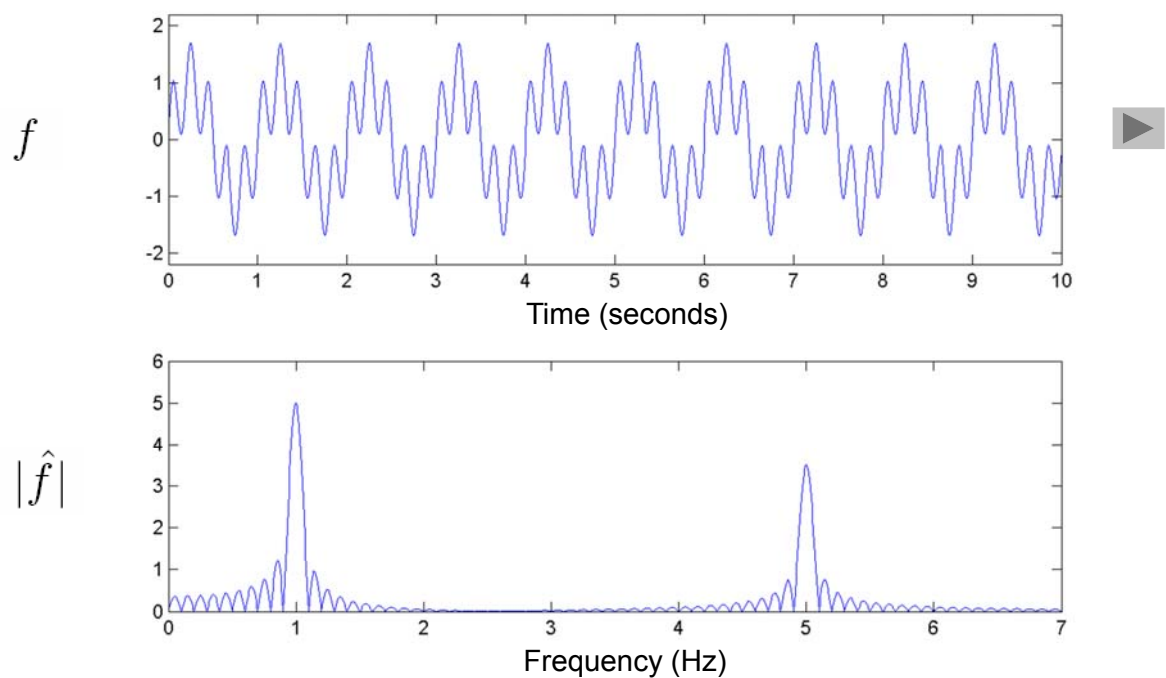Lecture

**Music Processing**

# Audio Features

**Meinard Müller**

International Audio Laboratories Erlangen
meinard.mueller@audiolabs-erlangen.de

---

## Fourier Transform

# Fourier Transform

Signal $\qquad\qquad\qquad\qquad f : \mathbb{R} \to \mathbb{R}$

Fourier representation $\quad f(t) = \int\limits_{\omega \in \mathbb{R}} c_\omega e^{2\pi i \omega t} d\omega \ , \ c_\omega = \hat{f}(\omega)$

Fourier transform $\qquad\quad \hat{f}(\omega) = \int\limits_{t \in \mathbb{R}} f(t) e^{-2\pi i \omega t} dt$
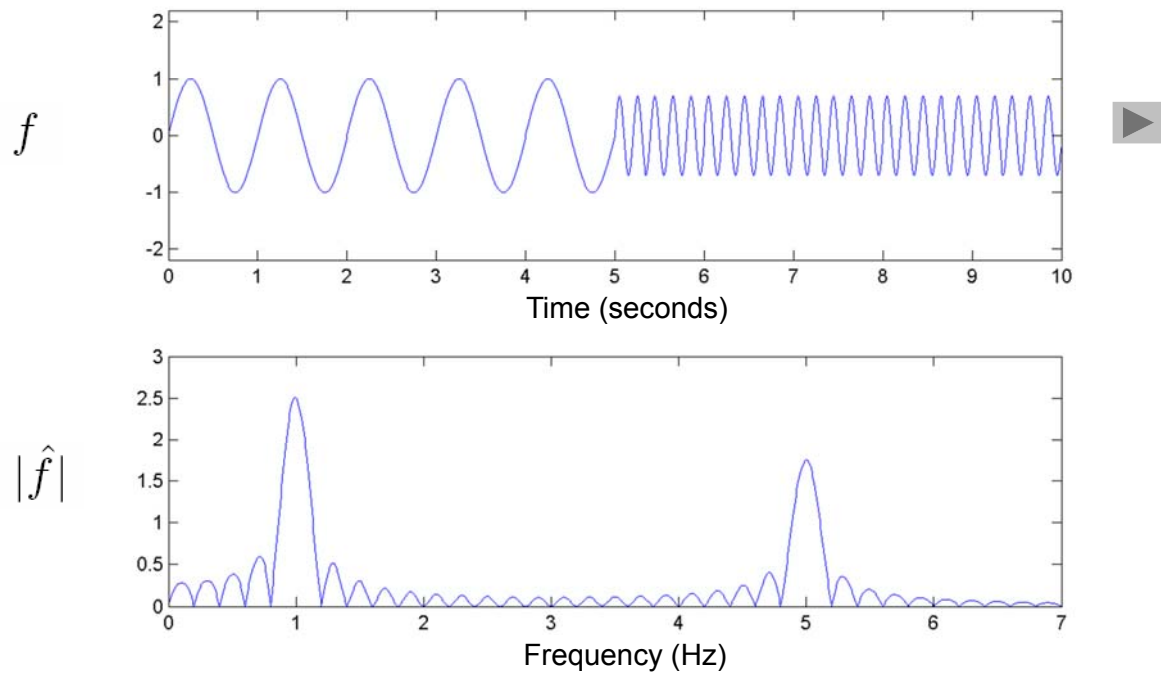
---

# Fourier Transform

Signal $\qquad\qquad\qquad\qquad f : \mathbb{R} \to \mathbb{R}$

Fourier representation $\quad f(t) = \int\limits_{\omega \in \mathbb{R}} c_\omega e^{2\pi i \omega t} d\omega \ , \ c_\omega = \hat{f}(\omega)$

Fourier transform $\qquad\quad \hat{f}(\omega) = \int\limits_{t \in \mathbb{R}} f(t) e^{-2\pi i \omega t} dt$

- 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
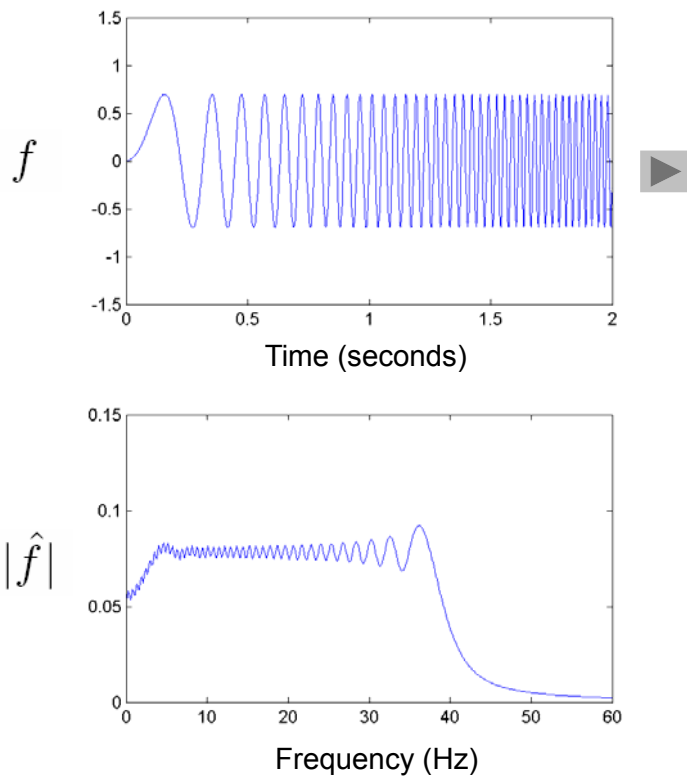
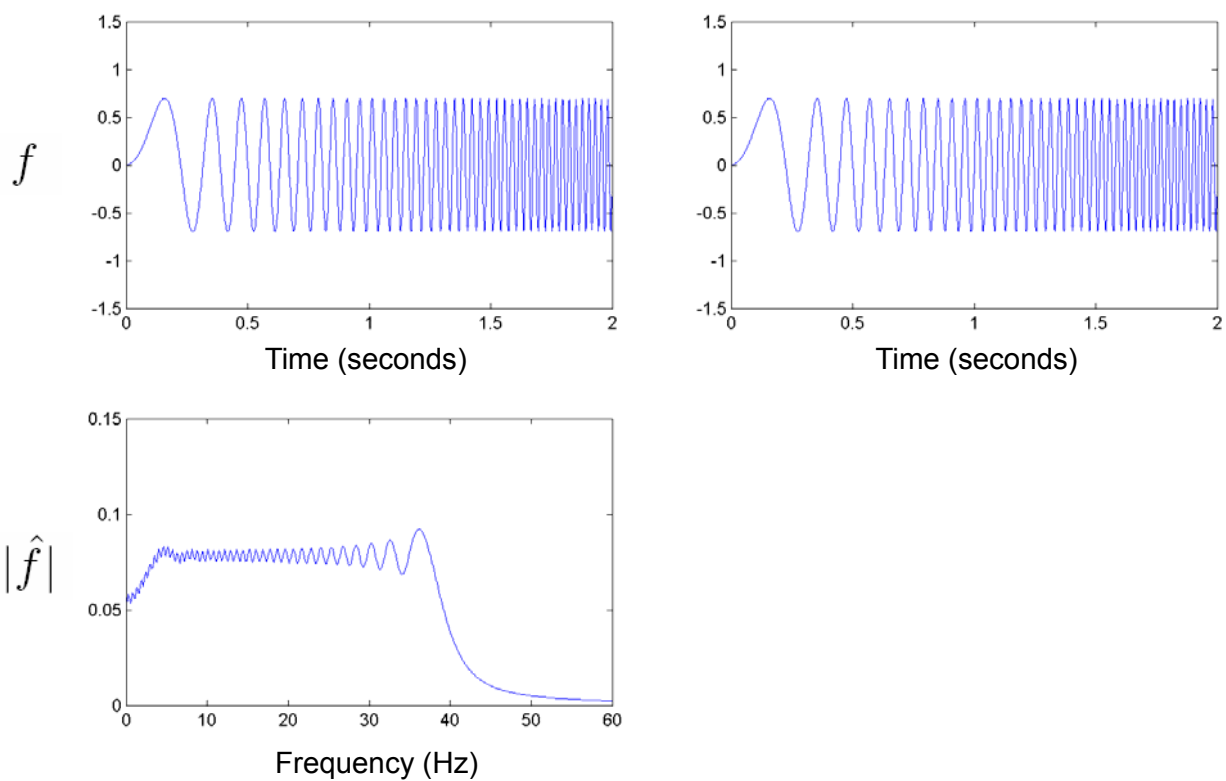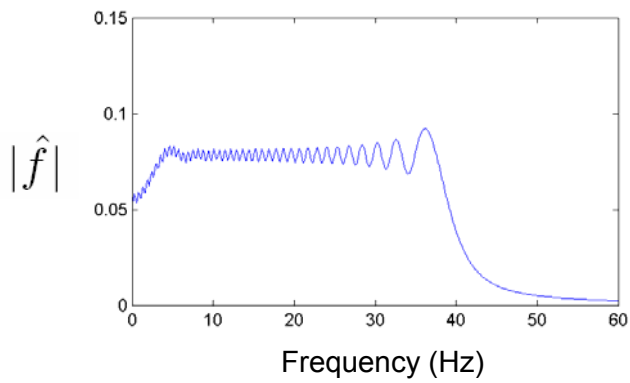# Fourier Transform



# Short Time Fourier Transform

Idea (Dennis Gabor, 1946):

- Consider only a small section of the signal
  for the spectral analysis

  $\rightarrow$ recovery of time information

- Short Time Fourier Transform (STFT)

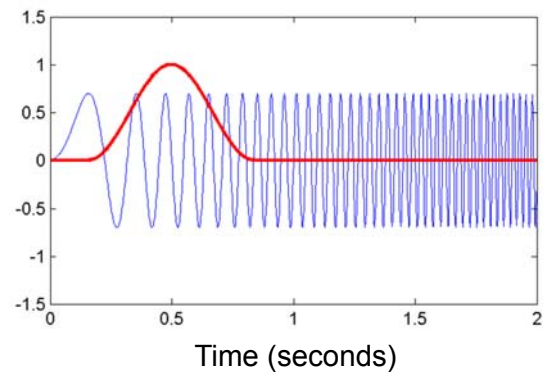- Section is determined by pointwise multiplication
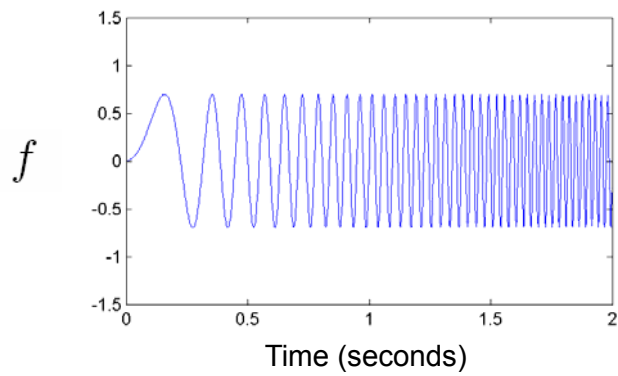  of the signal with a localizing window function

# Short Time Fourier Transform



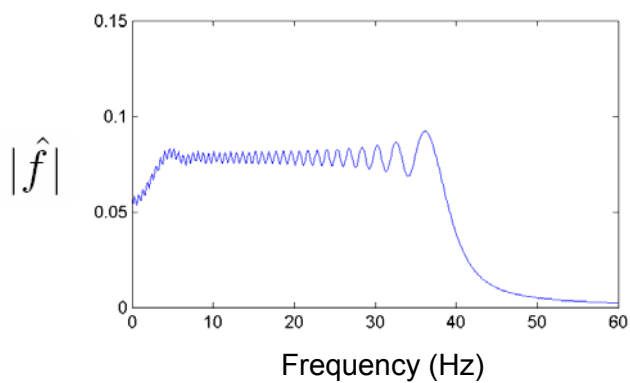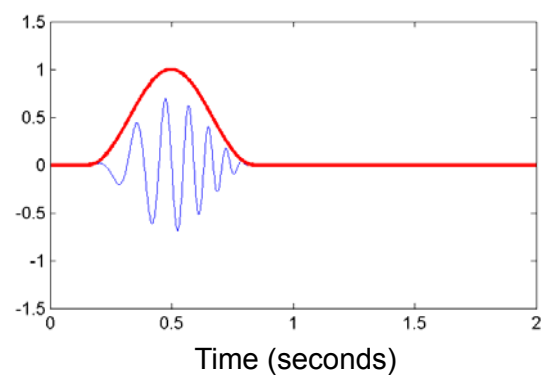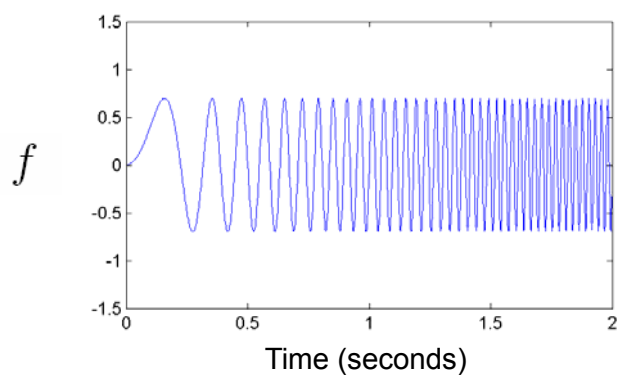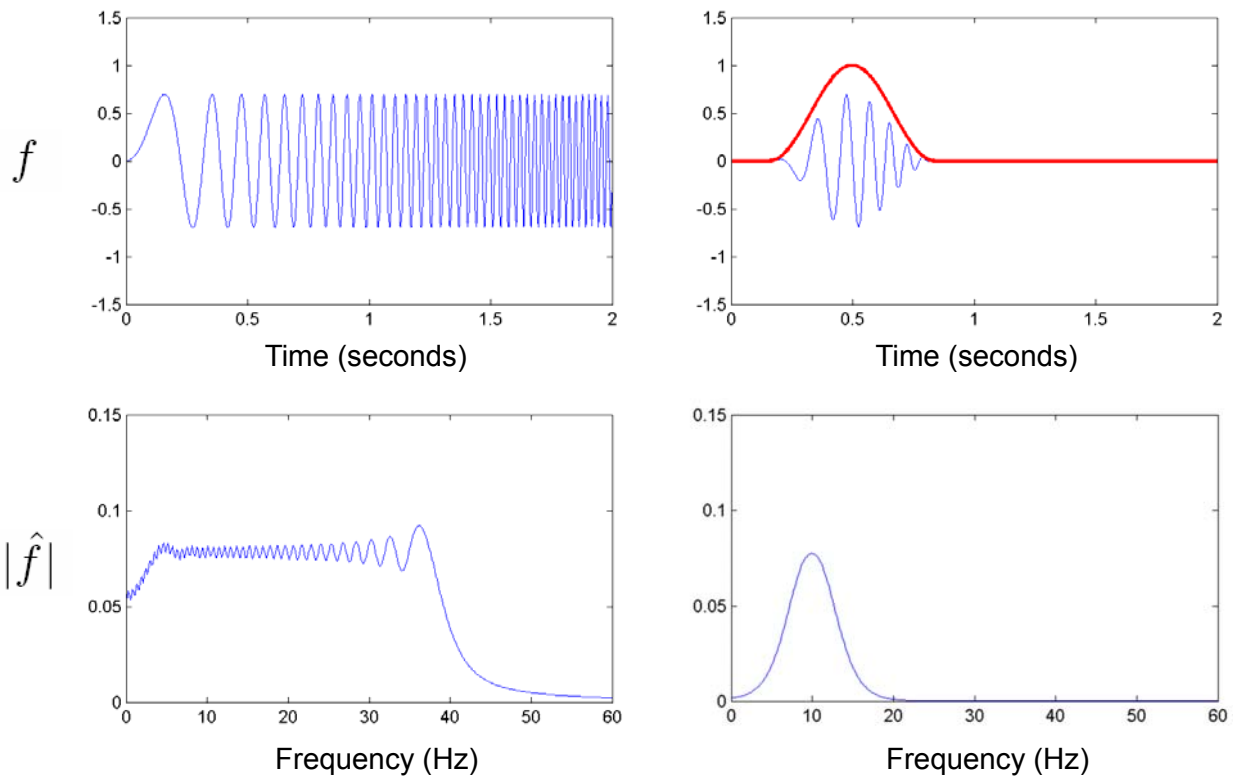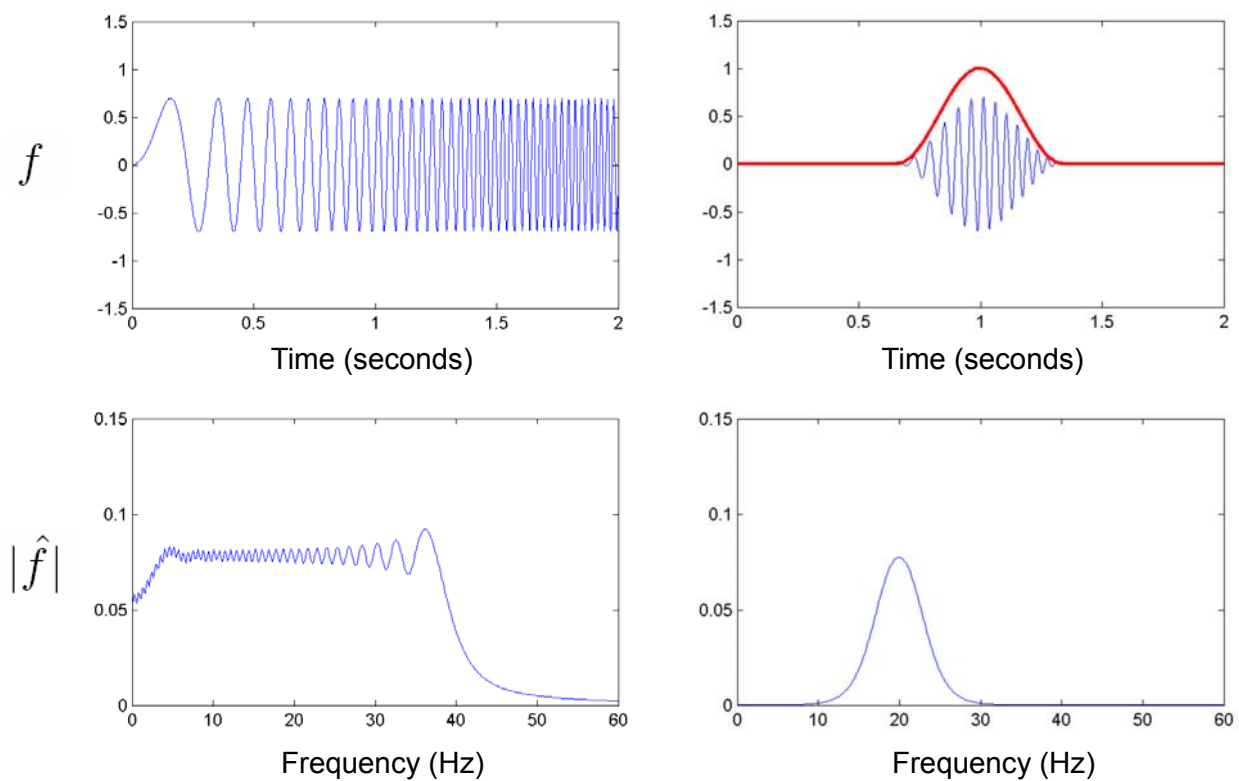# Short Time Fourier Transform

# Short Time Fourier Transform


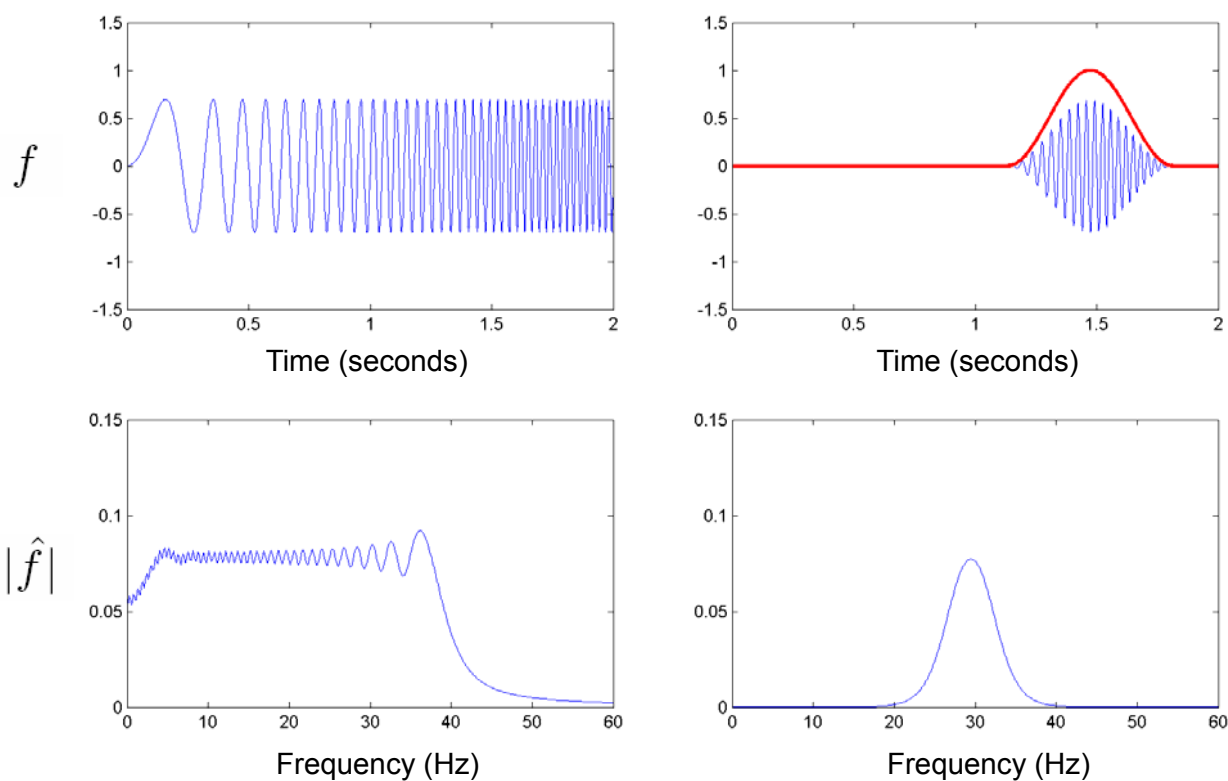
# Short Time Fourier Transform

# Short Time Fourier Transform



# Short Time Fourier Transform

# Short Time Fourier Transform



---

# Short Time Fourier Transform

Definition

- Signal $\qquad f : \mathbb{R} \to \mathbb{R}$

- Window function $\quad g : \mathbb{R} \to \mathbb{R} \quad (g \in L^2(\mathbb{R}), \|g\| = 1)$

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

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

# Short Time Fourier Transform

Intuition:

- $g_{\omega,t}$ is „musical note" of frequency $\omega$, which oscillates within the translated window $u \to g(u - t)$



# Short Time Fourier Transform

Intuition:

- $g_{\omega,t}$ is „musical note" of frequency $\omega$, which oscillates within the translated window $u \to g(u - t)$
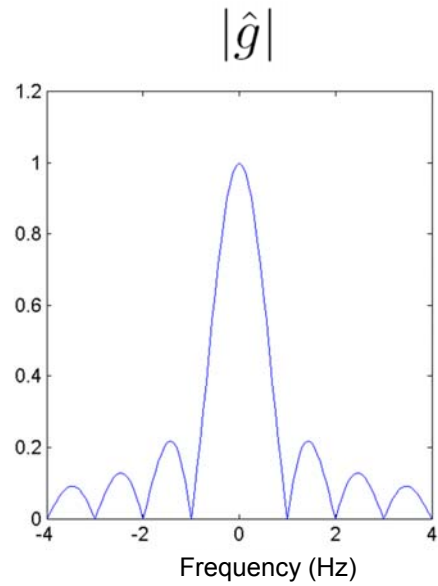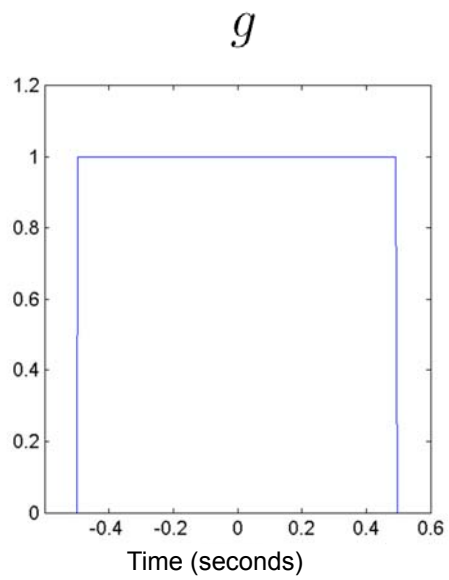


- Innere product $\langle f | g_{\omega,t} \rangle$ measures the correlation between the musical note $g_{\omega,t}$ and the signal $f$.
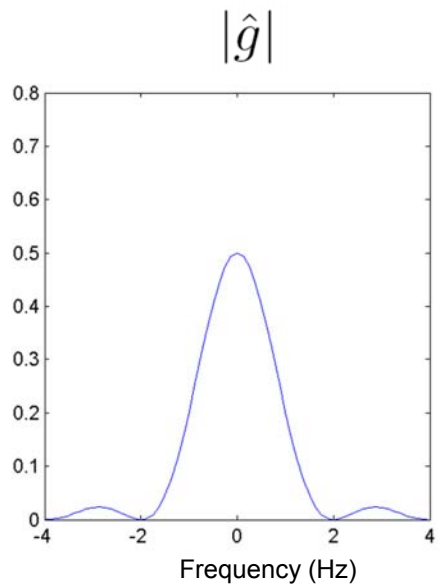
# Window Function

Box window
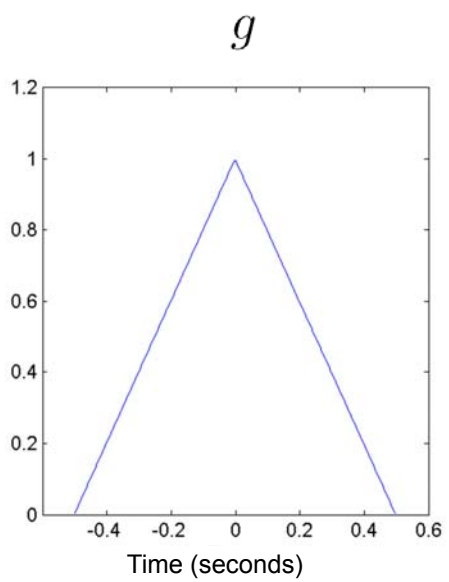

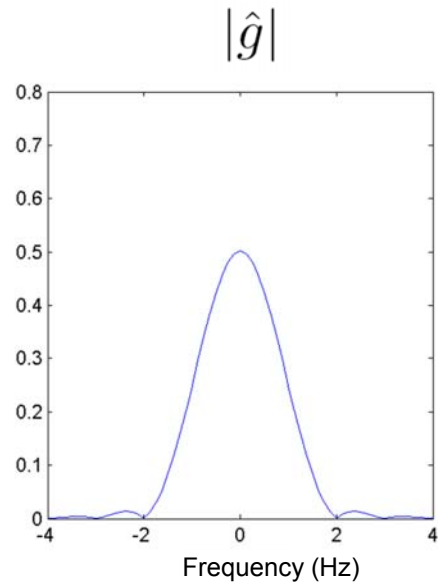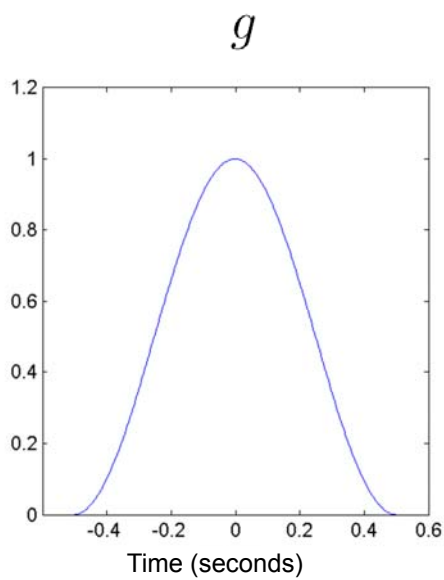
$g$ — Time (seconds)

$|\hat{g}|$ — Frequency (Hz)

---

# Window Function

Triangle window



$g$ — Time (seconds)

$|\hat{g}|$ — Frequency (Hz)
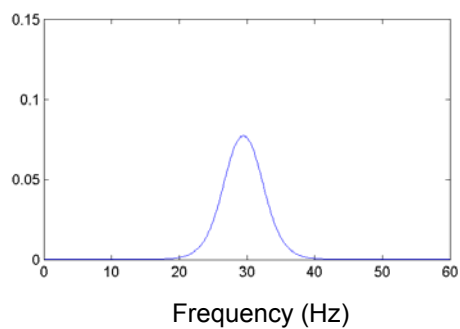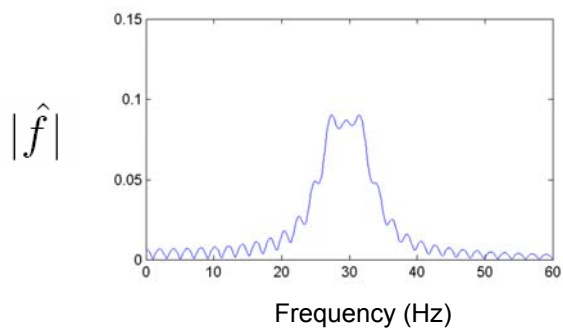
# Window Function

Hann window



# Window Function



Trade off between smoothing and „ringing"

# Time-Frequency Representation

$f$

Frequency $\omega$
(Hertz)

$|\langle f | g_{\omega,t} \rangle|$

Time $t$ (seconds)

Intensity
(dB)

# Time-Frequency Representation

$f$

Frequency $\omega$
(Hertz)

$|\langle f | g_{\omega,t} \rangle|$

**Spectrogram**

Time $t$ (seconds)

Intensity
(dB)

# Time-Frequency Representation

Chirp signal and STFT with box window of length 0.05



# Time-Frequency Representation

Chirp signal and STFT with Hann window of length 0.05

# Time-Frequency Localization

- Size of window constitutes a trade-off 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.

---

# Short Time Fourier Transform

Signal and STFT with Hann window of length 0.02

# Short Time Fourier Transform

Signal and STFT with Hann window of <span style="color:red">length 0.1</span>



# Heisenberg Uncertainty Principle

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

| Center | Width |
|---|---|

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

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

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

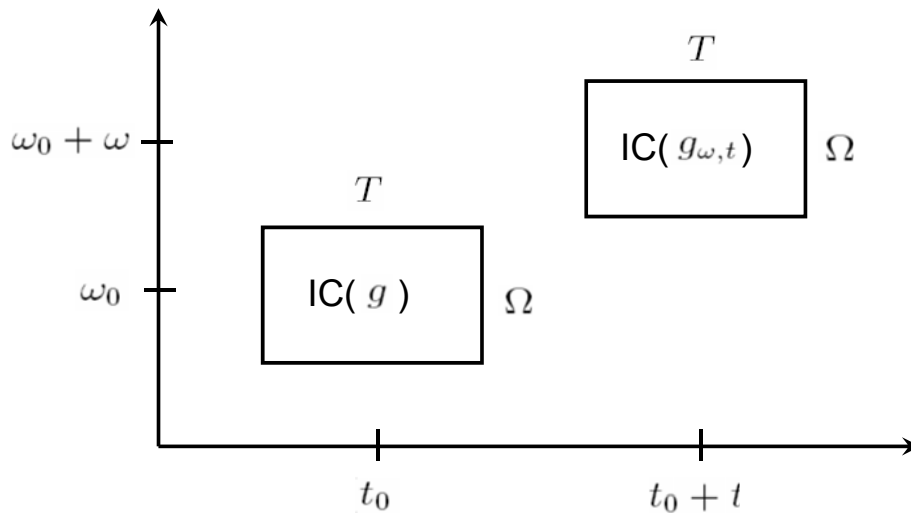# Information Cells

$$g_{\omega,t}(u) := e^{2\pi i \omega u} g(u - t) \quad \text{``musical note''}$$



# MATLAB

- MATLAB function SPECTROGRAM
- $N =$ window length (in samples)
- $M =$ overlap (usually $N/2$)
- Compute $\text{DFT}_N$ for every windowed section
- Keep lower $N/2$ Fourier coefficients

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

# Example

Let $x$ be a discrete time signal $\quad x(n) = f(Tn)$

Sampling rate: $\quad 1/T = 22050 \text{ Hz}$

Window length: $\quad N = 4096$

Overlap: $\quad N/2 = 2048$

Hopsize: $\quad \text{window length} - \text{overlap}$

$$
\begin{aligned}
\text{Let} \quad v_0 \; &:= \; (x(0),\; x(1),\; \ldots,\; x(4095)) \\
v_1 \; &:= \; (x(2048),\; \ldots,\; x(6143)) \\
v_2 \; &:= \; (x(4096),\; \ldots,\; x(8191)) \\
&\;\vdots
\end{aligned}
$$

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

---

# Example

Time resolution:

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

Frequency resolution:

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

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

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

# Pitch Features

Model assumption:    Equal-tempered scale

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

$\rightarrow$ Logarithmic frequency distribution

Octave: doubling of frequency

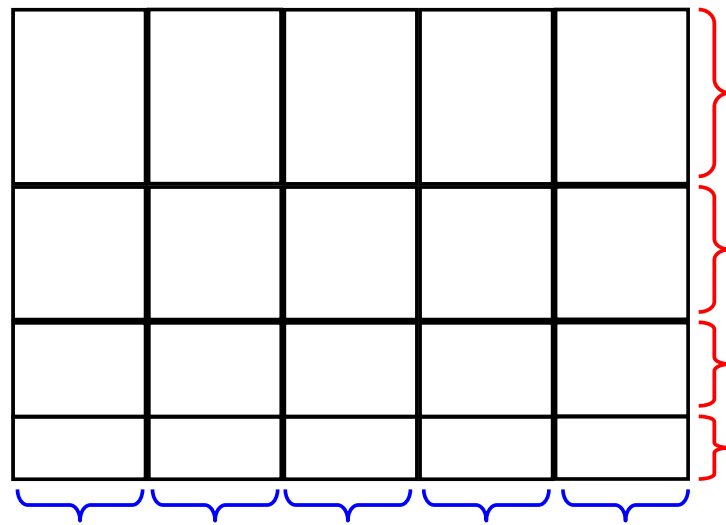# Pitch Features

Idea: Binning of Fourier coefficients

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

## Pitch Features

Time-frequency representation



Windowing in the time domain

Windowing in the frequency domain

## 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

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

# Pitch Features

Example: A4, *p* = 69

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

$$
\begin{aligned}
f(k = 79) &= 425.3 \; Hz \\
f(k = 80) &= 430.7 \; Hz \\
f(k = 81) &= 436.0 \; Hz \\
f(k = 82) &= 441.4 \; Hz \\
f(k = 83) &= 446.8 \; Hz \\
f(k = 84) &= 452.2 \; Hz \\
f(k = 85) &= 457.6 \; Hz
\end{aligned}
$$

---

# Pitch Features

Example: A4, *p* = 69

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

$$
\begin{aligned}
f(k = 79) &= 425.3 \; Hz \\
f(k = 80) &= 430.7 \; Hz \\
f(k = 81) &= 436.0 \; Hz \\
f(k = 82) &= 441.4 \; Hz \\
f(k = 83) &= 446.8 \; Hz \\
f(k = 84) &= 452.2 \; Hz \\
f(k = 85) &= 457.6 \; Hz
\end{aligned}
$$

$S(p = 69)$

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

## Pitch Features

| Note | MIDI pitch | Center [Hz] frequency | 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 |

## 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.

$\rightarrow$ Linear frequency sampling is problematic!

Solution:

Multi-resolution spectrograms or multirate filterbanks

# Pitch Features

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



# Pitch Features



Spectrogram

# Pitch Features

## Spectrogram



# Pitch Features

## Pitch representation
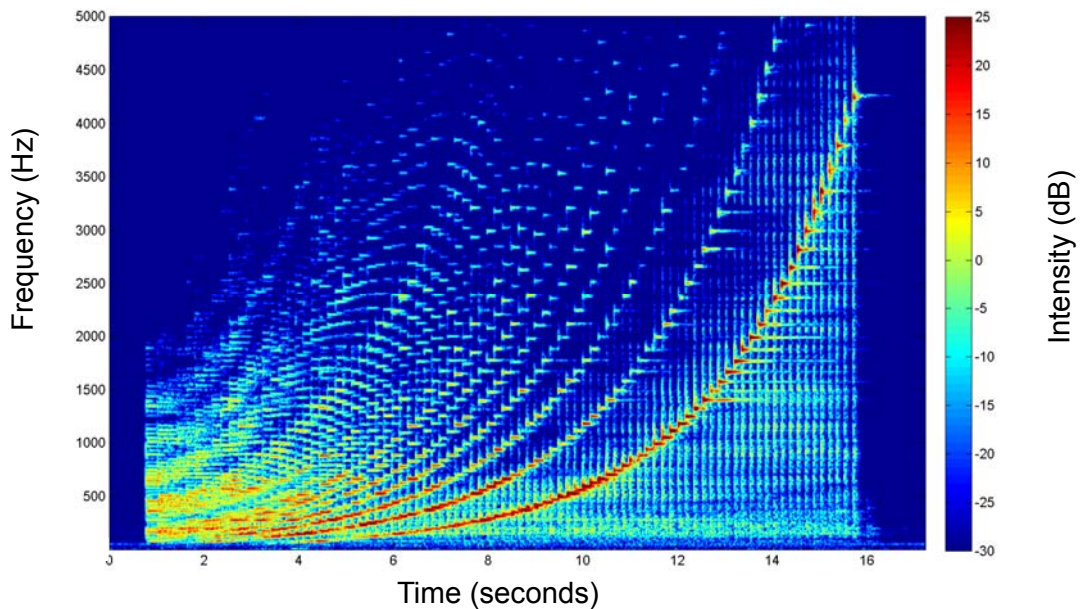
# Pitch Features



## Pitch representation



---

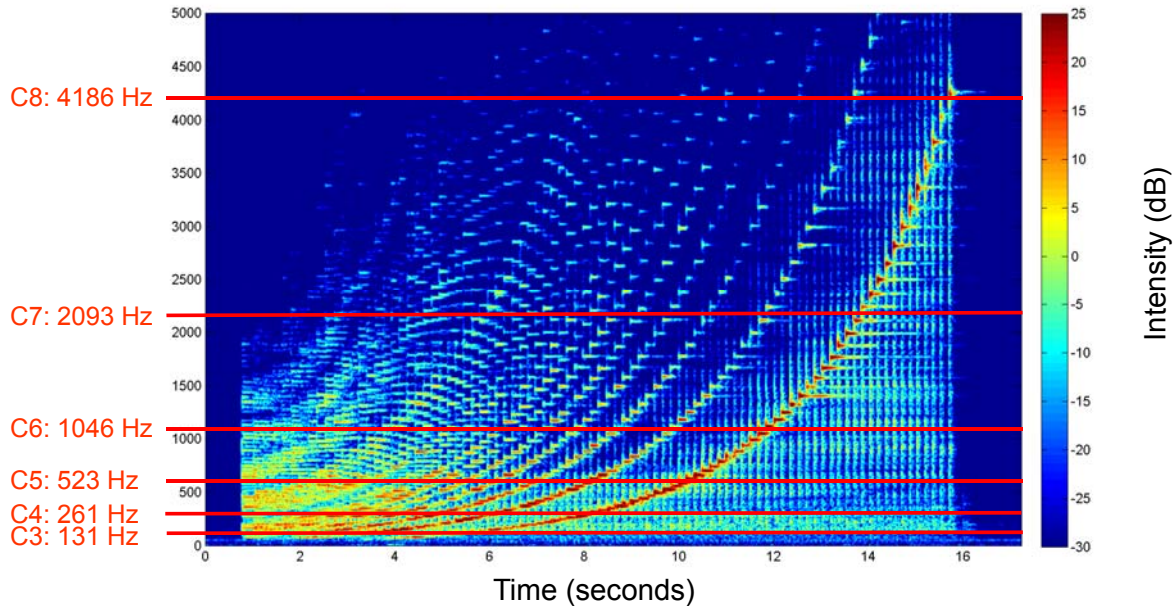# Pitch Features

Example: Chromatic scale

## Spectrogram
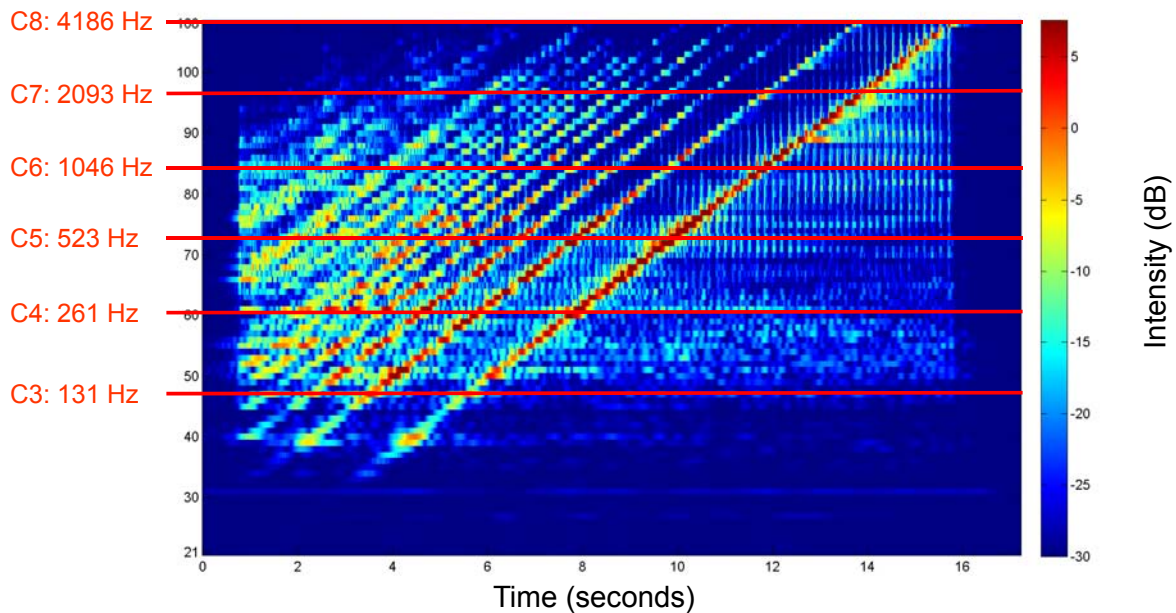
# Pitch Features

Example: Chromatic scale

Spectrogram



# Pitch Features
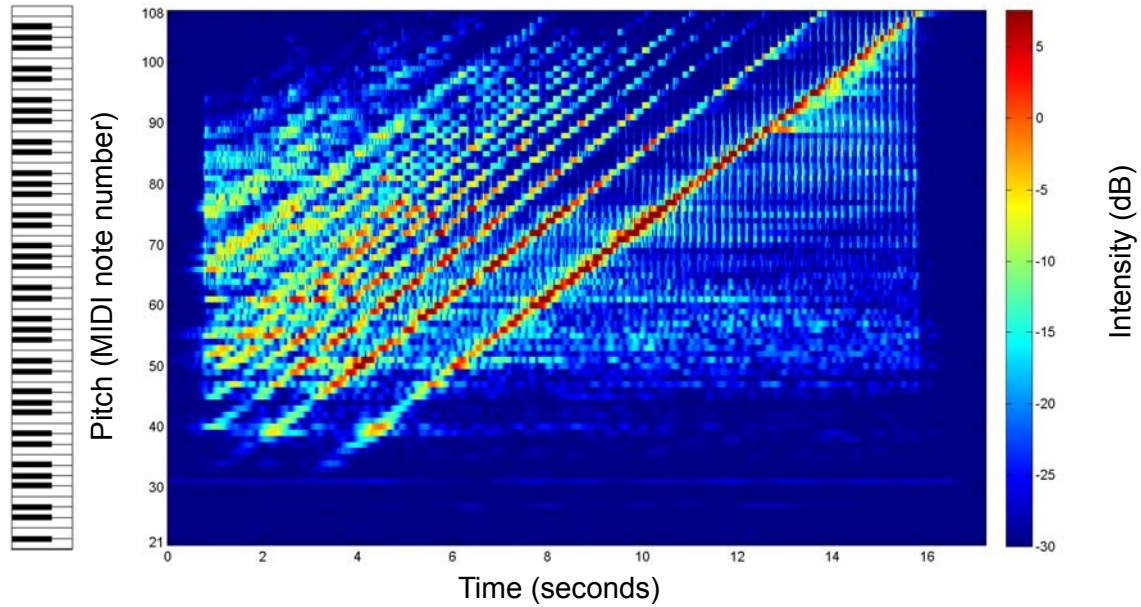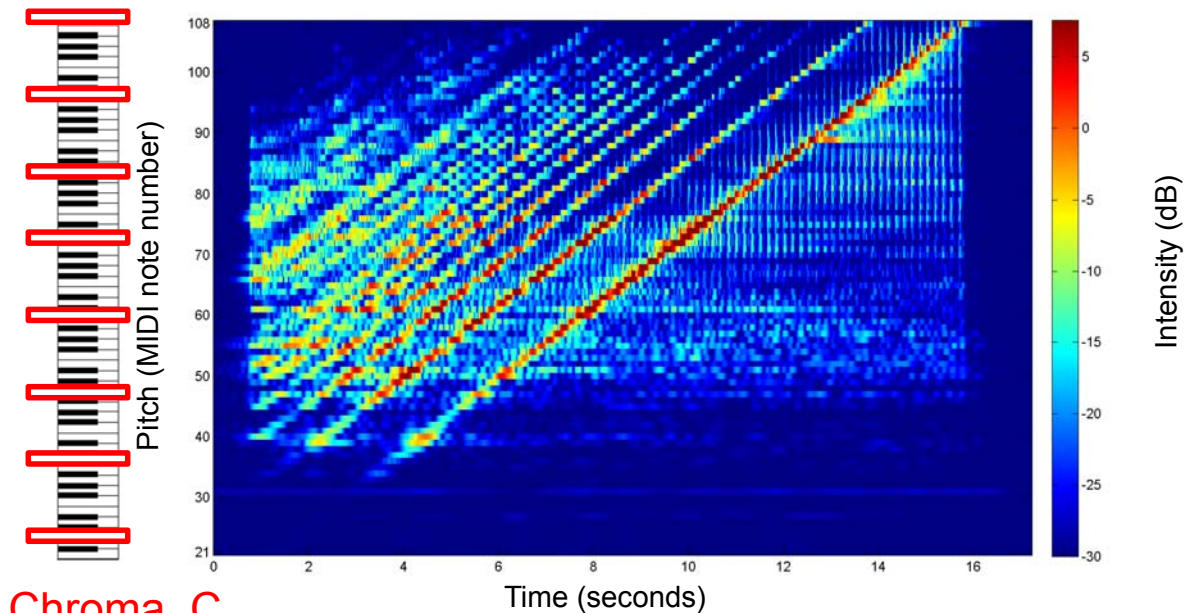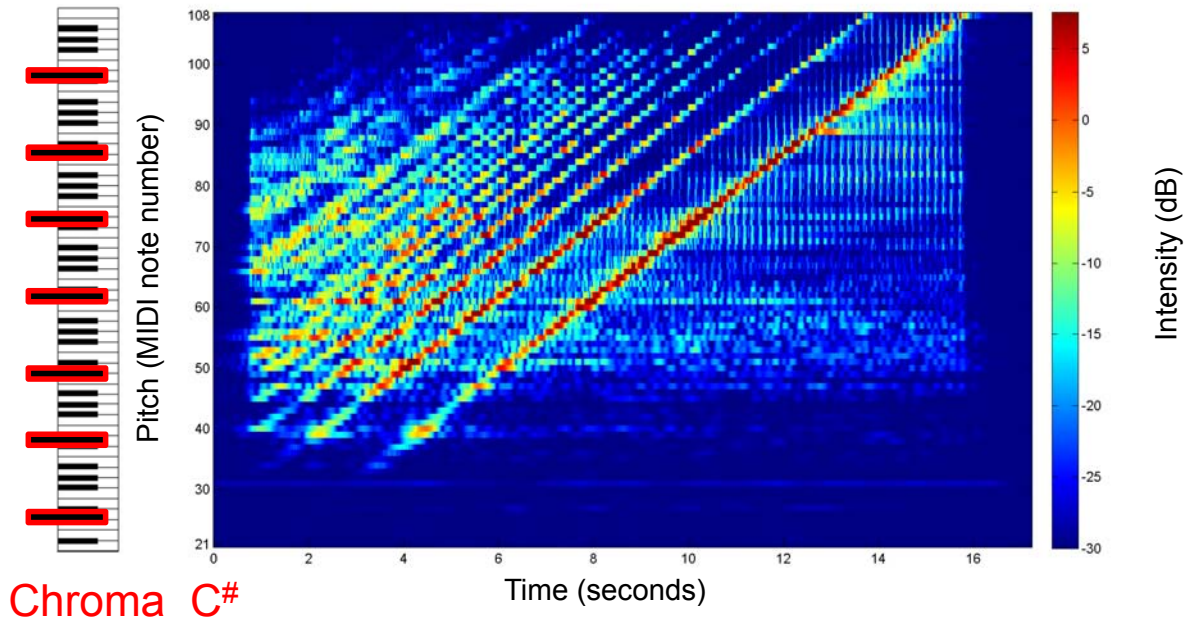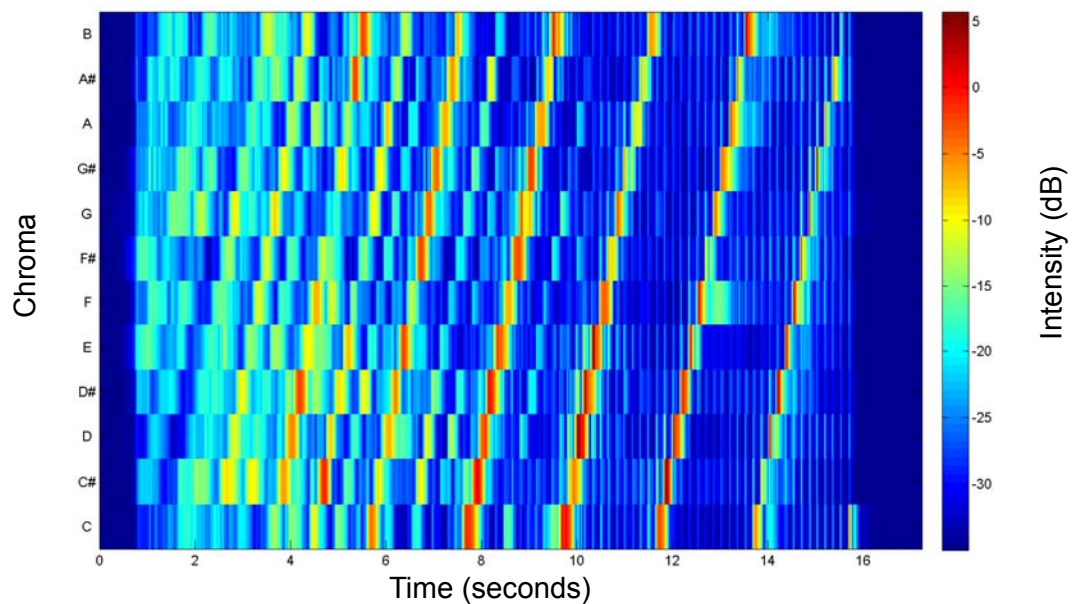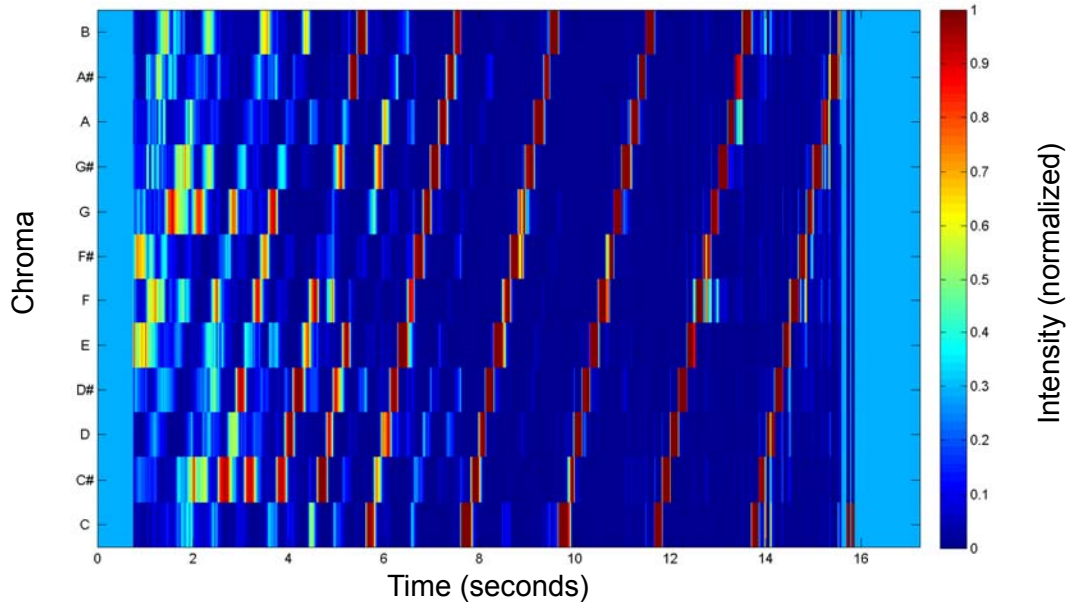
Example: Chromatic scale

Log-frequency spectrogram

# Pitch Features

Example: Chromatic scale

Log-frequency spectrogram



---

# Pitch Features

Example: Chromatic scale

Log-frequency spectrogram



Chroma  C

# Pitch Features

Example: Chromatic scale ▶

**Log-frequency spectrogram**



**Chroma  C#**

# Chroma Features

Example: Chromatic scale ▶

**Chroma representation** ▶

# Chroma Features

Example: Chromatic scale

<span style="color:red">Chroma representation (normalized, Euclidean)</span>



---

# 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.
- Seperation of pitch into two components: <span style="color:red">tone height</span> (octave number) and <span style="color:red">chroma</span>.
- Chroma : 12 traditional pitch classes of the equal-tempered scale. For example:

  Chroma C $\; \widehat{=} \; \{\ldots, \; C0, \; C1, \; C2, \; C3, \; \ldots\}$
- Computation: pitch features $\rightarrow$ chroma features

  Add up all pitches belonging to the same class
- Result: 12-dimensional chroma vector.

# Chroma Features



# Chroma Features



C2          C3          C4

Chroma C

# Chroma Features



C#2        C#3        C#4
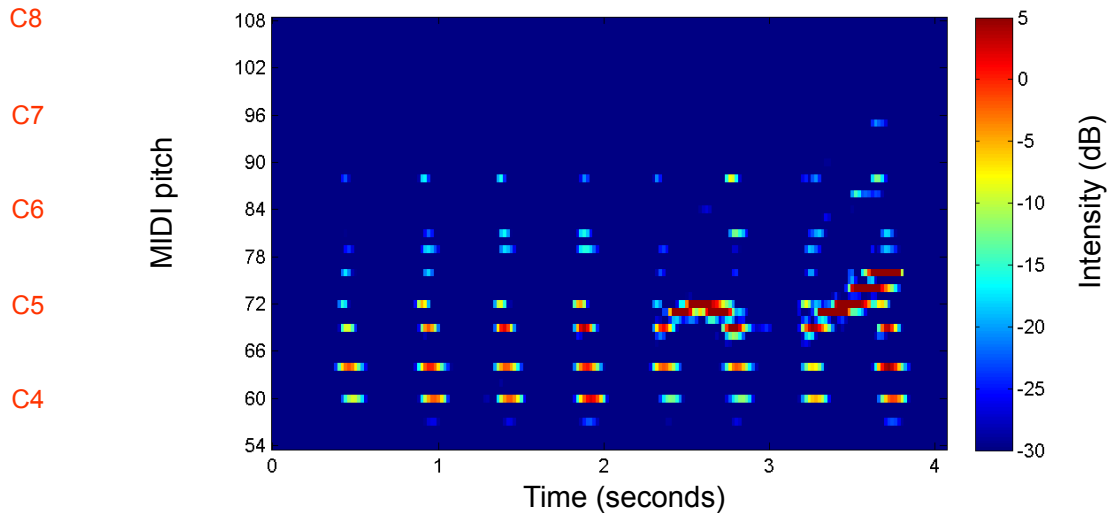
Chroma  C#

# Chroma Features



D2        D3        D4
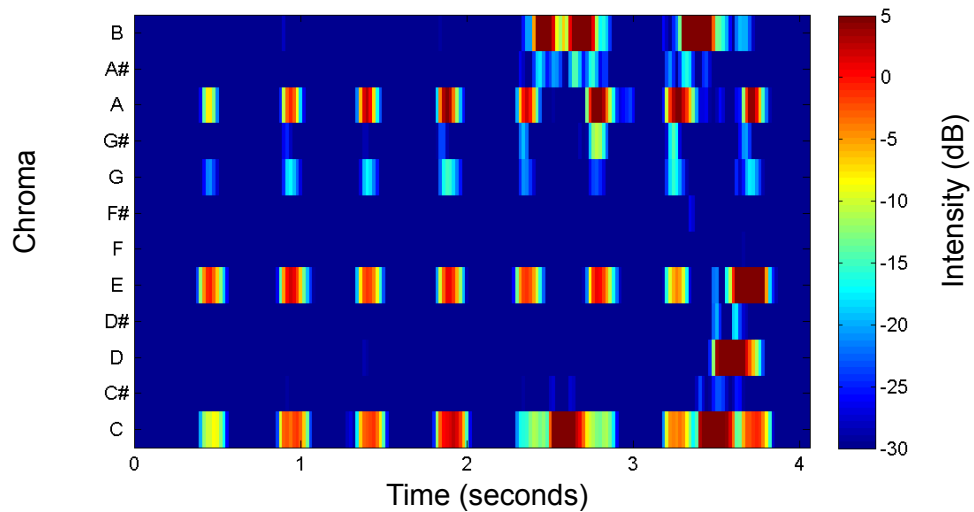
Chroma  D

# Chroma Features

Chromatic circle        Shepard's helix of pitch perception



http://en.wikipedia.org/wiki/Pitch_class_space        Bartsch/Wakefield, IEEE Trans. Multimedia, 2005

---

# Chroma Features

Example: C-Major Scale        ▶        ▶

# Chroma Features



## Pitch representation



---

# Chroma Features



## Chroma representation

# Chroma Features



Chroma representation (normalized)

# Chroma Features

Example: Beethoven's Fifth
Chroma representation (normalized, 10 Hz)

Karajan  ▶

Scherbakov  ▶

# Chroma Features

Example: Beethoven's Fifth
Chroma representation (normalized, 2 Hz)
Smoothing (2 seconds) + downsampling (factor 5)



# Chroma Features
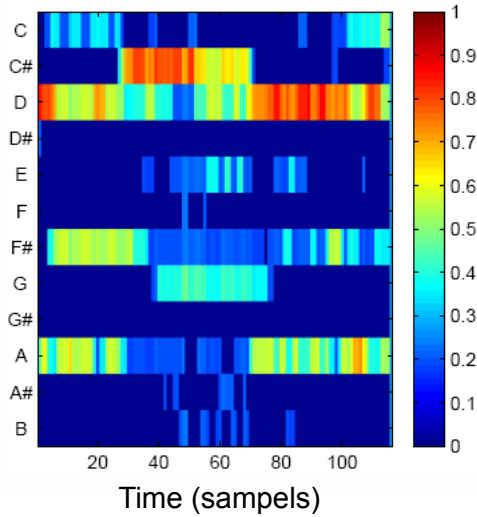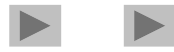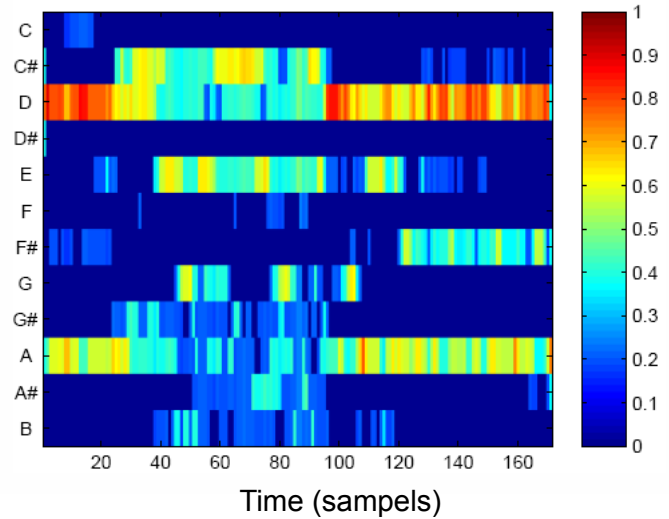
Example: Beethoven's Fifth
Chroma representation (normalized, 1 Hz)
Smoothing (4 seconds) + downsampling (factor 10)

# Chroma Features

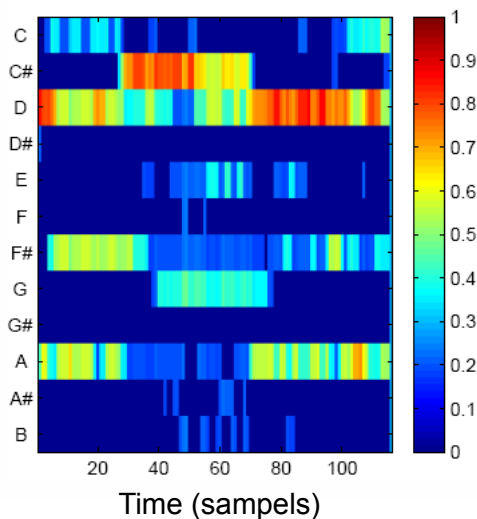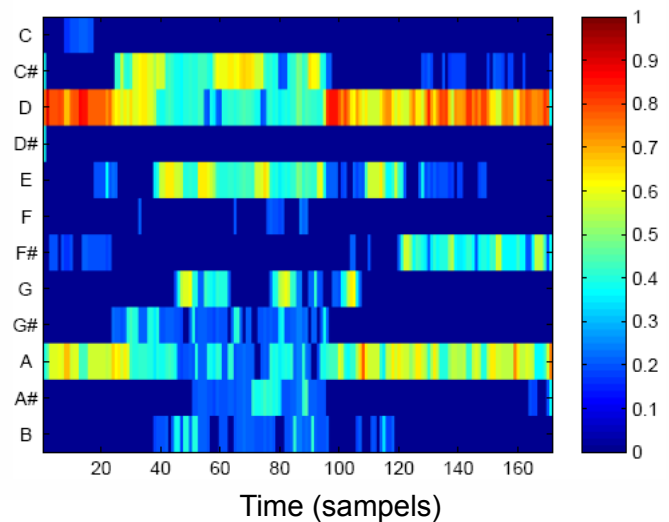Example: Bach Toccata

Koopman ▶ ▶   Ruebsam ▶ ▶



# Chroma Features

Example: Bach Toccata
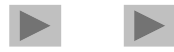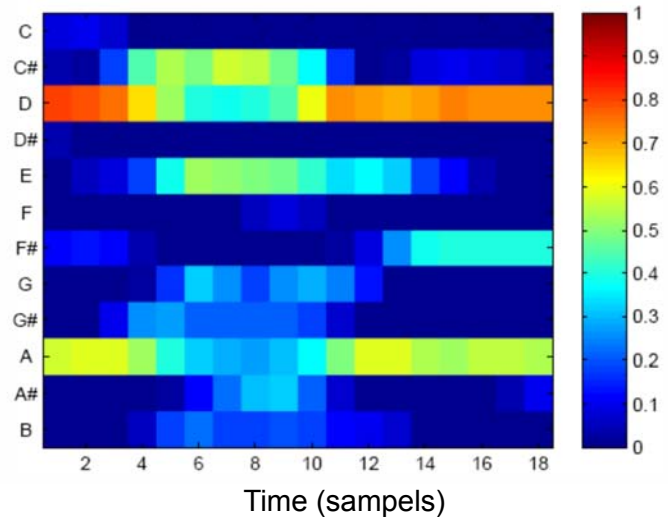
Koopman ▶ ▶   Ruebsam ▶ ▶


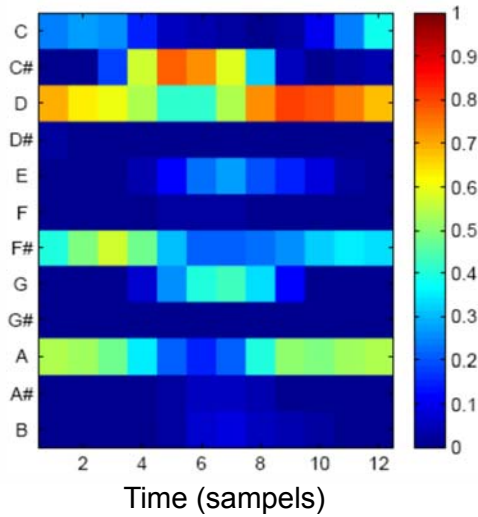
Feature resolution: 10 Hz

# Chroma Features

Example: Bach Toccata

Koopman  ▶  ▶          Ruebsam  ▶  ▶

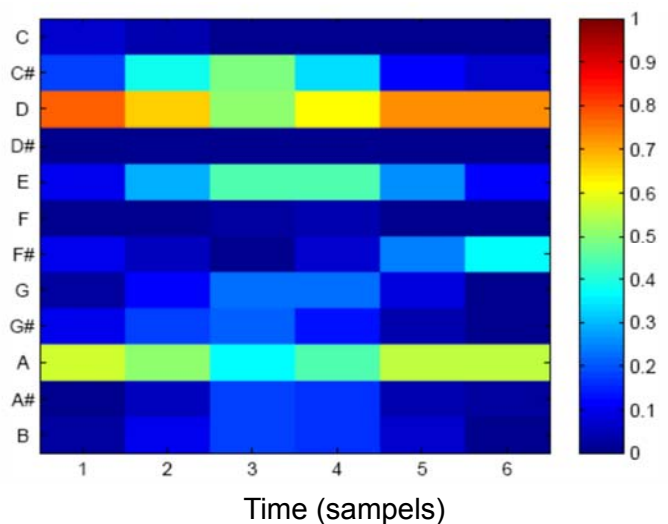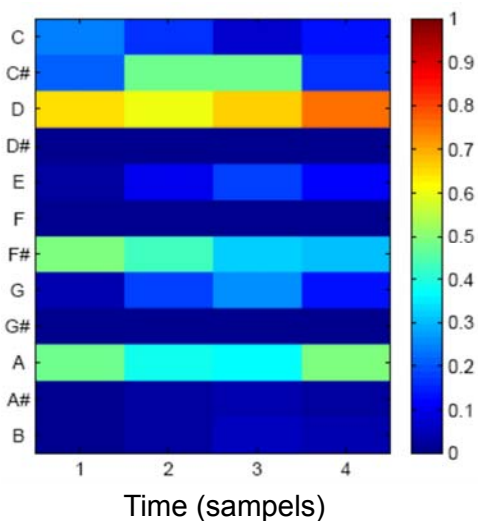Feature resolution: 1 Hz

# Chroma Features
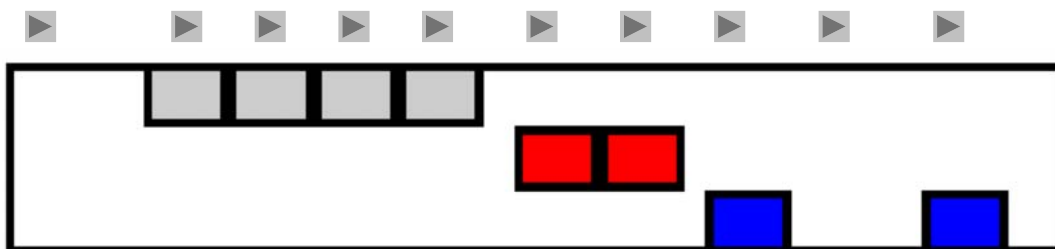
Example: Bach Toccata

Koopman  ▶  ▶          Ruebsam  ▶  ▶

Feature resolution: 0.33 Hz

# Chroma Features

- Sequence of chroma vectors correlates to the harmonic progression

- Normalization $v \to \dfrac{v}{\|v\|}$ 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
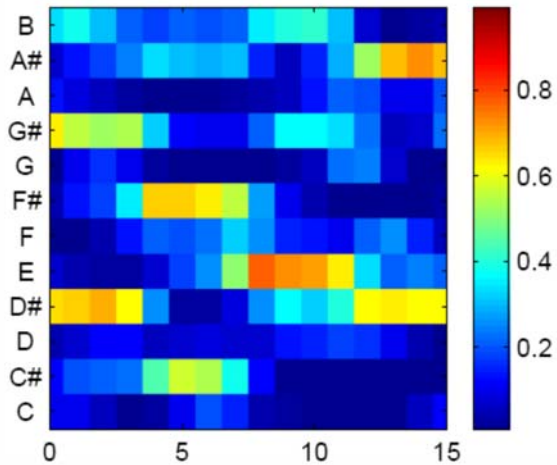
---

# Chroma Features

Example: Zager & Evans "In The Year 2525"
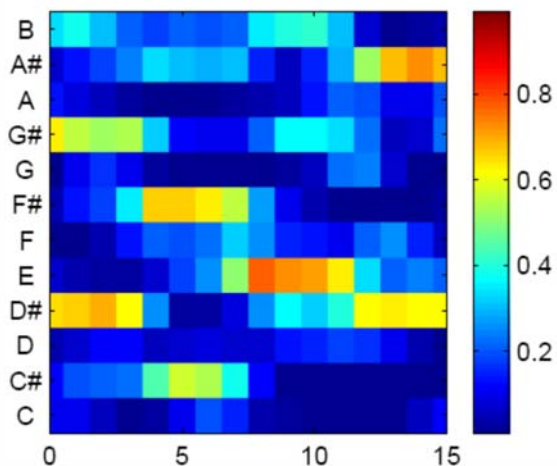


How to deal with transpositions?

# Chroma Features

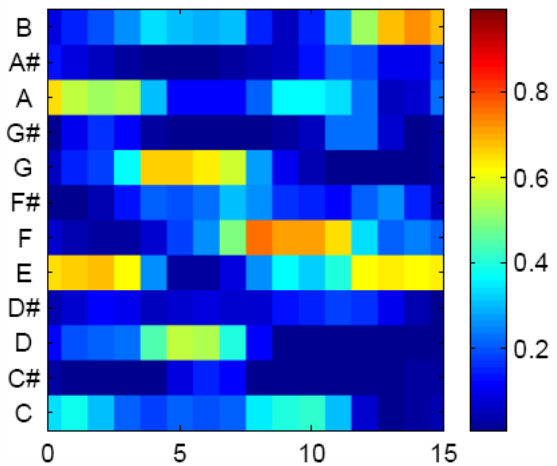Example: Zager & Evans "In The Year 2525"



Original: $(v^1, \ldots, v^N)$

---

# Chroma Features

Example: Zager & Evans "In The Year 2525"
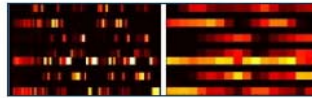


Original: $(v^1, \ldots, v^N)$



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

# Audio Features

- There are many ways to implement chroma features
- Properties may differ significantly
- Appropriateness depends on respective application



Chroma Toolbox: Pitch, Chroma, CENS, CRP

- http://www.mpi-inf.mpg.de/resources/MIR/chromatoolbox/
- MATLAB implementations for various chroma variants
- ISMIR 2011, Poster Session (PS2), Tuesday 13-15