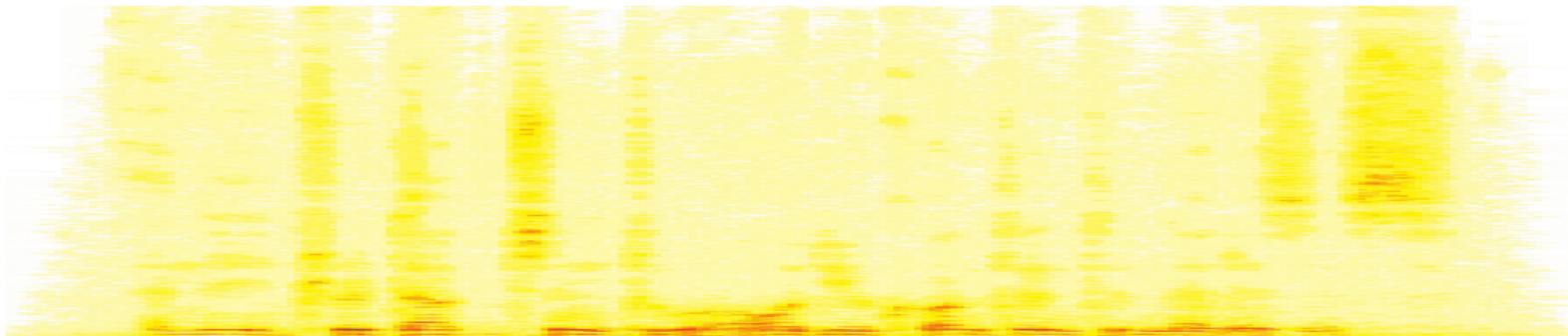


# Introduction to Audio Content Analysis

## Module 5.7: Musical Key Recognition

alexander lerch



# introduction

## overview

corresponding textbook section

[Chapter 5 — Tonal Analysis](#): pp. 88–94

[Chapter 5 — Tonal Analysis](#): pp. 116–125

### ● lecture content

- definition of musical key
- pitch chroma feature
- standard approach for key recognition

### ● learning objectives

- explain the defining properties of a musical key
- implement a simple pitch chroma feature extractor
- describe and discuss a simple automatic key recognition system



# introduction

## overview

corresponding textbook section

[Chapter 5 — Tonal Analysis](#): pp. 88–94

[Chapter 5 — Tonal Analysis](#): pp. 116–125

### ● lecture content

- definition of musical key
- pitch chroma feature
- standard approach for key recognition

### ● learning objectives

- explain the defining properties of a musical key
- implement a simple pitch chroma feature extractor
- describe and discuss a simple automatic key recognition system



# key

## tonic & mode

- **tonic:** first scale degree
  - most “important” pitch class
  
- **mode:** set of diatonic pitch relationships
  - Major: 2, 2, 1, 2, 2, 2, 1
  - Minor: 2, 1, 2, 2, 1, 2, 2

Major

(Aeolic) Minor

(Harmonic) Minor

Dorian

Phrygian

Lydian

Mixolydian

Lokrian

Chromatic

Whole tone

# key

## key & key signature 1/2

- **key:**

defined by *tonic* (root note) and *mode*

- defines a set of pitch classes constructing both pitch and harmonic content

- **modulation** (local key changes):

common in various styles, uncommon in others

- **key signature:**

indicates current key with accidentals (score notation)

# key

## key & key signature 1/2

- **key:**

defined by *tonic* (root note) and *mode*

- defines a set of pitch classes constructing both pitch and harmonic content

- **modulation** (local key changes):

common in various styles, uncommon in others

- **key signature:**

indicates current key with accidentals (score notation)

# key

## key & key signature 1/2

- **key:**

defined by *tonic* (root note) and *mode*

- defines a set of pitch classes constructing both pitch and harmonic content

- **modulation** (local key changes):

common in various styles, uncommon in others

- **key signature:**

indicates current key with accidentals (score notation)

# key

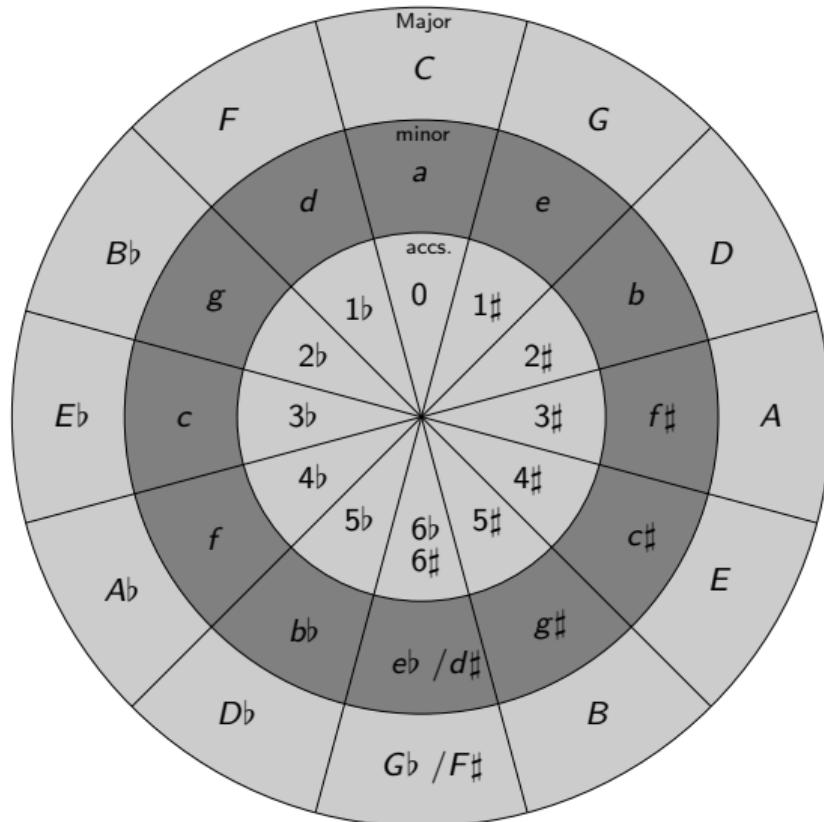
## key & key signature 2/2

The image displays a vertical stack of eight musical staves, each consisting of five horizontal lines. The staves are separated by vertical bar lines. Below each staff, the name of a major key is written in capital letters. The keys are arranged in two columns of four. The first column contains C Major, G Major, D Major, A Major, E Major, B Major, F# Major, and Gb Major. The second column contains Db Major, Ab Major, Eb Major, Bb Major, F Major, and C Major. The notes in each staff are represented by small circles (heads) on the lines, indicating the pitch. The first staff (C Major) has no sharps or flats. Subsequent staves introduce sharps (#) or flats (b) in a consistent pattern across the two columns.

C Major      G Major  
D Major      A Major  
E Major      B Major  
F<sup>#</sup> Major      G<sup>b</sup> Major  
Db Major      Ab Major  
Eb Major      Bb Major  
F Major      C Major

# musical pitch

key: circle of fifths

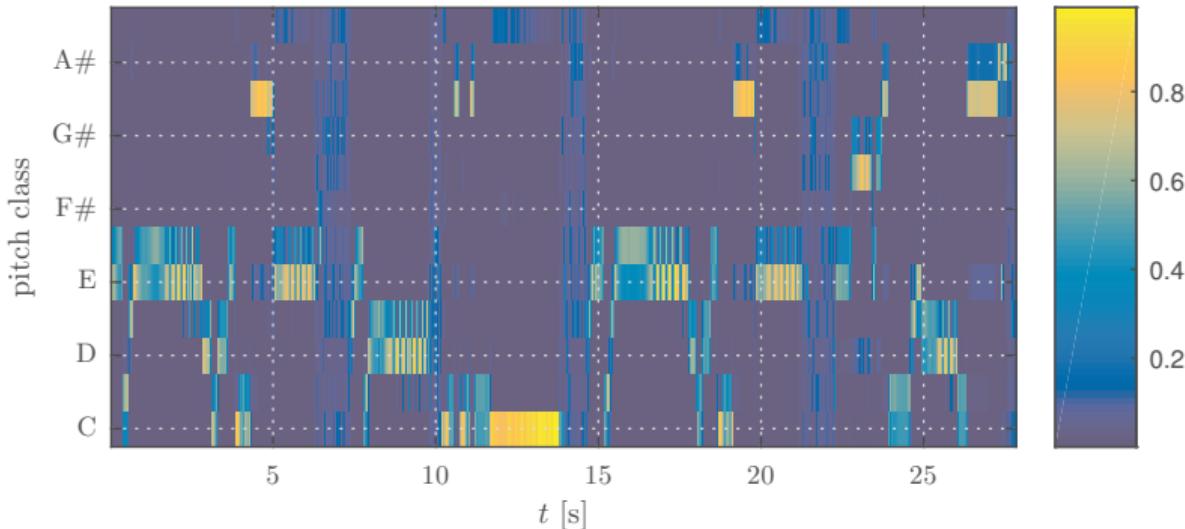


# pitch chroma

## introduction



- pitch class distribution
- 12-dimensional vector

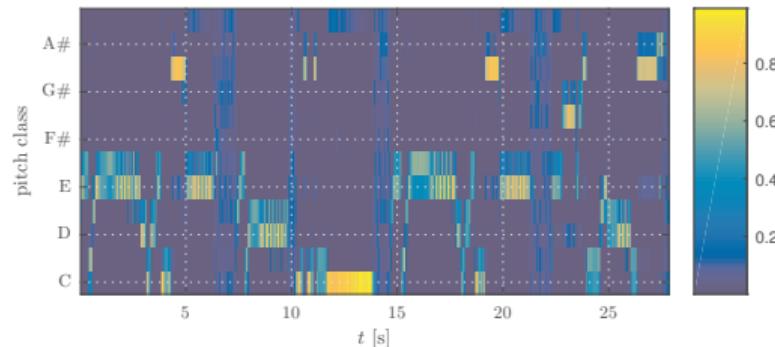


# pitch chroma

## introduction



- pitch class distribution
- 12-dimensional vector



- no octave information
  - robust representation
  - no differentiation between prime and octave

# pitch chroma computation 1/2

- ➊ divide spectral representation into **semi-tone bands**
- ➋ compute **mean** per band

$$\mu(j, n) = \frac{1}{k_u(j) - k_l(j) + 1} \sum_{k=k_l(j)}^{k_u(j)} |X(k, n)|$$

- ➌ sum/mean every 12th band

$$\begin{aligned}\nu(j\%12, n) &= \sum_{o=o_l}^{o_u} \mu(j, n), \\ \nu(n) &= [\nu(0, n), \nu(1, n), \nu(2, n), \dots, \nu(10, n), \nu(11, n)]^T\end{aligned}$$

# pitch chroma computation 1/2

- ➊ divide spectral representation into **semi-tone bands**
- ➋ compute **mean per band**

$$\mu(j, n) = \frac{1}{k_u(j) - k_l(j) + 1} \sum_{k=k_l(j)}^{k_u(j)} |X(k, n)|$$

- ➌ sum/mean every 12th band

$$\begin{aligned}\nu(j\%12, n) &= \sum_{o=o_l}^{o_u} \mu(j, n), \\ \nu(n) &= [\nu(0, n), \nu(1, n), \nu(2, n), \dots, \nu(10, n), \nu(11, n)]^T\end{aligned}$$

# pitch chroma computation 1/2

- ➊ divide spectral representation into **semi-tone bands**
- ➋ compute **mean per band**

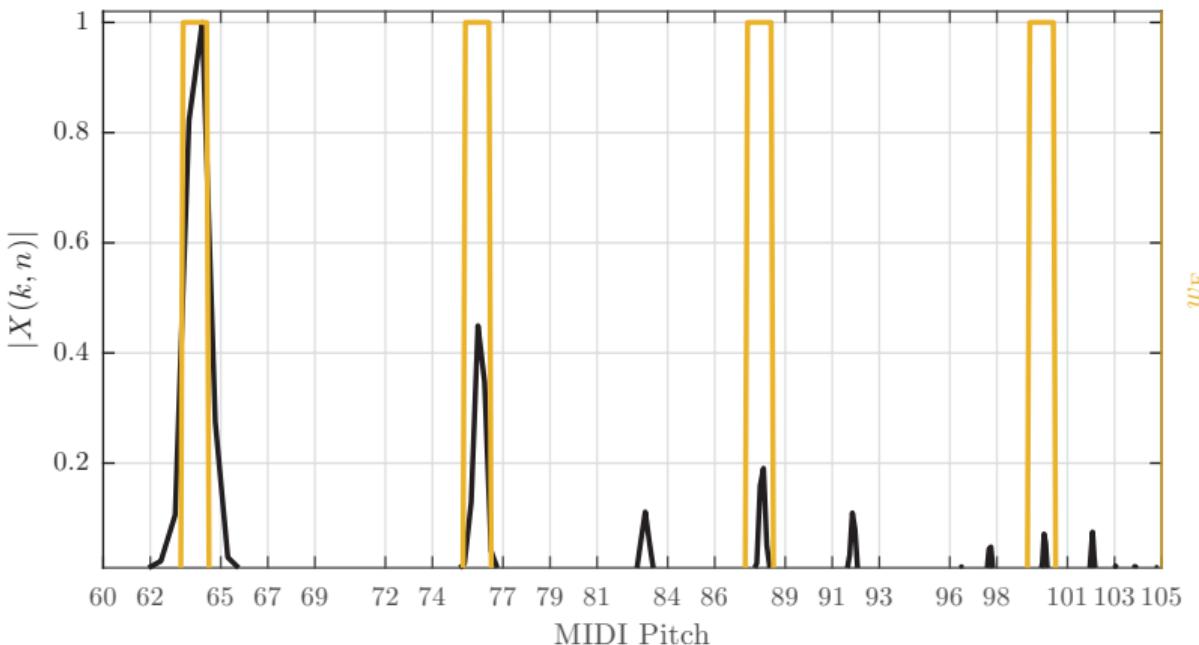
$$\mu(j, n) = \frac{1}{k_u(j) - k_l(j) + 1} \sum_{k=k_l(j)}^{k_u(j)} |X(k, n)|$$

- ➌ sum/mean every 12th band

$$\begin{aligned}\nu(j\%12, n) &= \sum_{o=o_l}^{o_u} \mu(j, n), \\ \nu(n) &= [\nu(0, n), \nu(1, n), \nu(2, n), \dots, \nu(10, n), \nu(11, n)]^T\end{aligned}$$

# pitch chroma

## computation 2/2

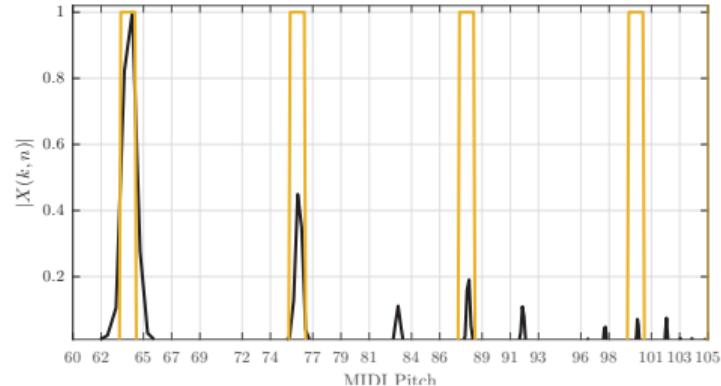


# pitch chroma

computation: simple variants

- **STFT:**

- *weighted mean of bins (window function)*
- *tonalness preprocessing (local maxima etc)*



- sum of filterbank output energies

- **CQT:**

- sum of bins/peaks

- beat-synchronous chroma

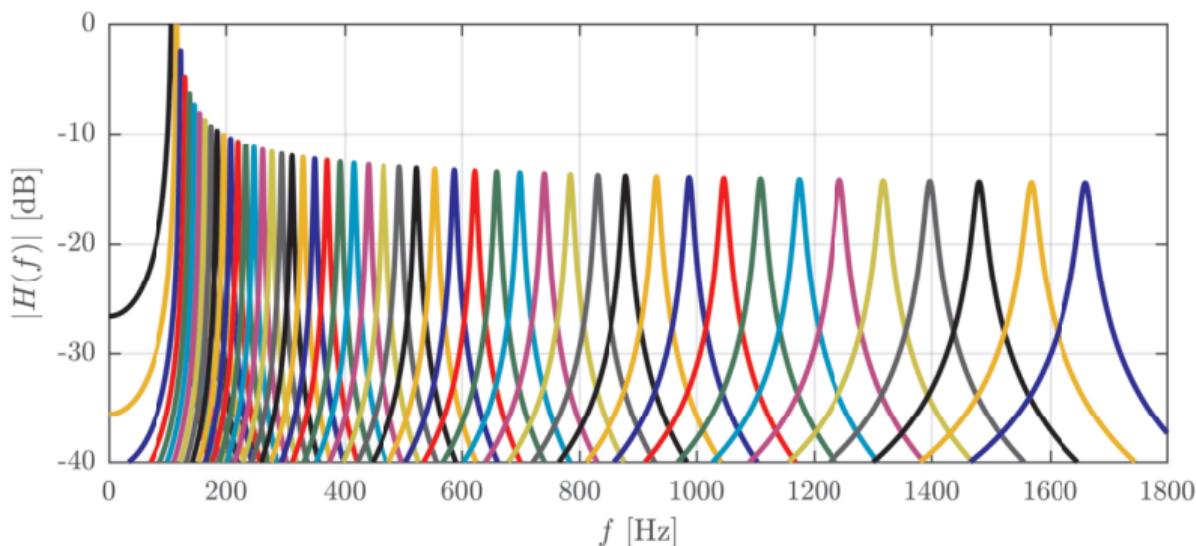
# pitch chroma

## computation: simple variants

- **STFT:**

- weighted mean of bins (window function)
- tonalness preprocessing (local maxima etc)

- sum of **filterbank** output energies



# pitch chroma

computation: simple variants

- **STFT:**

- *weighted mean of bins (window function)*
- *tonalness preprocessing (local maxima etc)*

- sum of **filterbank** output energies

- **CQT:**

- sum of bins/peaks

- beat-synchronous chroma

# pitch chroma

computation: simple variants

- **STFT:**

- *weighted mean of bins (window function)*
- *tonalness preprocessing (local maxima etc)*

- sum of **filterbank** output energies

- **CQT:**

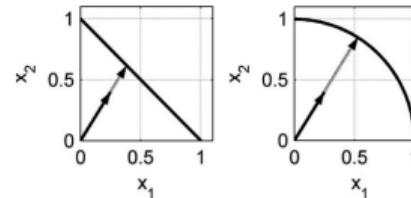
- sum of bins/peaks

- beat-synchronous chroma

# pitch chroma normalization

- pitch chroma as *distribution*:

$$\sum_{k=0}^{11} \nu(k, n) = 1$$



- pitch chroma as *vector*:

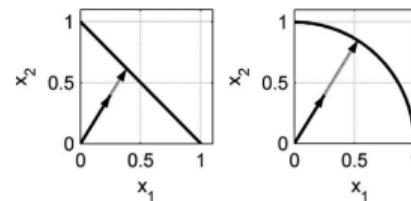
$$\sqrt{\sum_{k=0}^{11} \nu(k, n)^2} = 1$$

- other options:
  - e.g., short-term energy normalization (CENS)

# pitch chroma normalization

- pitch chroma as *distribution*:

$$\sum_{k=0}^{11} \nu(k, n) = 1$$



- pitch chroma as *vector*:

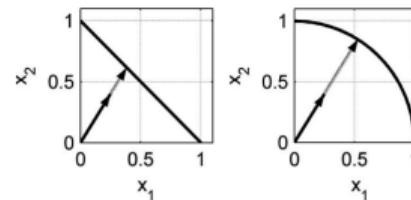
$$\sqrt{\sum_{k=0}^{11} \nu(k, n)^2} = 1$$

- other options:
  - e.g., short-term energy normalization (CENS)

# pitch chroma normalization

- pitch chroma as *distribution*:

$$\sum_{k=0}^{11} \nu(k, n) = 1$$



- pitch chroma as *vector*:

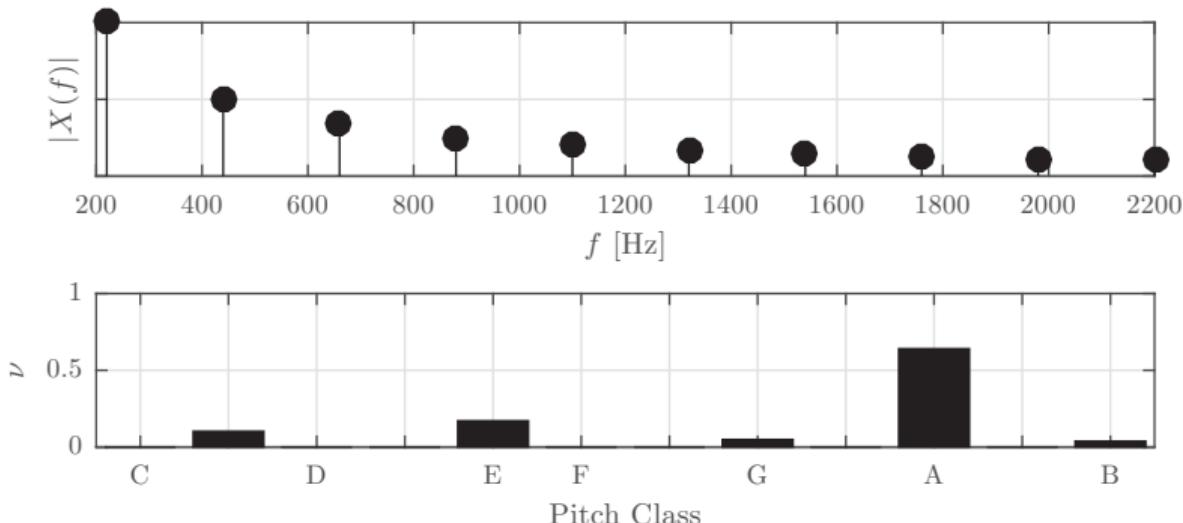
$$\sqrt{\sum_{k=0}^{11} \nu(k, n)^2} = 1$$

- other options:

- e.g., short-term energy normalization (CENS)

# pitch chroma

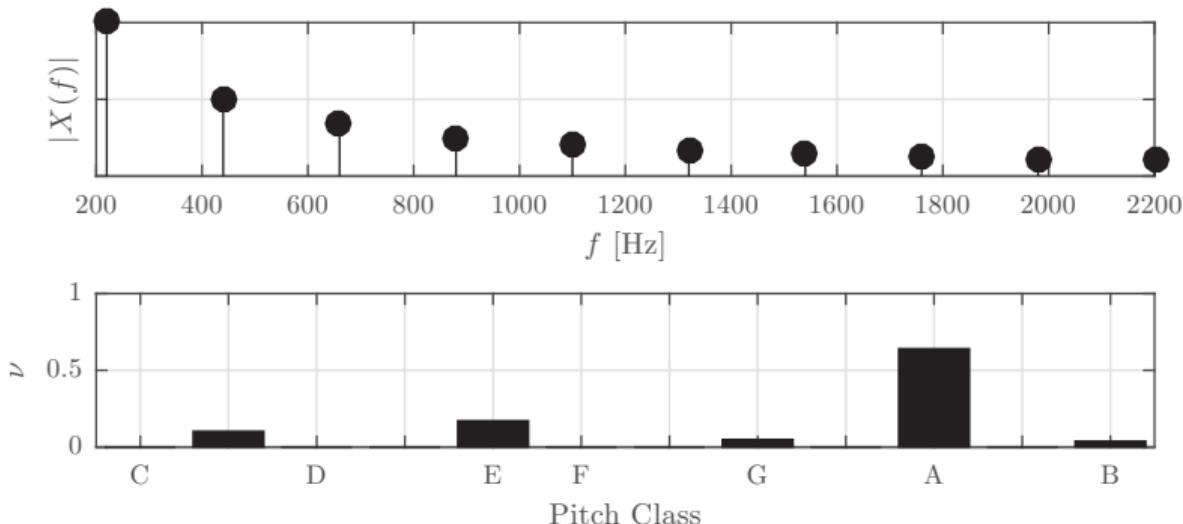
## problem 1: amplitude distortion



- every pitch contains not only fundamental but higher harmonics
  - de-emphasize higher frequencies
  - build amplitude model
  - use multi-pitch detection system

# pitch chroma

## problem 1: amplitude distortion



- every pitch contains not only fundamental but higher harmonics
  - ⇒ ⇒ de-emphasize higher frequencies
  - ⇒ build amplitude model
  - ⇒ use multi-pitch detection system

# pitch chroma

## problem 2: frequency distortion

- higher harmonics are not “in-tune”

Harmonic	$ \Delta C(f, f_T) $
$f = f_0$	0
$f = 2 \cdot f_0$	0
$f = 3 \cdot f_0$	1.955
$f = 4 \cdot f_0$	0
$f = 5 \cdot f_0$	13.6863
$f = 6 \cdot f_0$	1.955
$f = 7 \cdot f_0$	31.1741
$\mu_{ \Delta C }$	6.9672

# key detection

## introduction

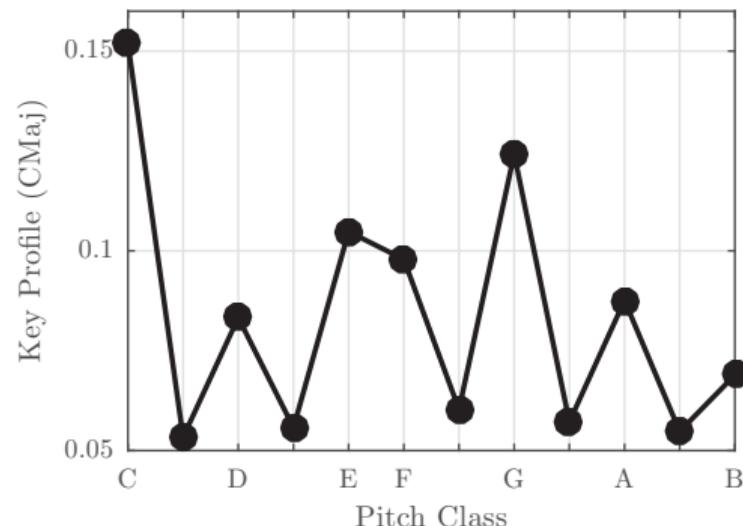
assumption:

- *pitch class distribution* is prototypical for key
  - tonic/root note is tonal center
  - tonal and harmonic relations define importance and occurrence of individual pitch classes
  - different root notes result in simple shift of distribution

# key detection

processing steps of simple key detection

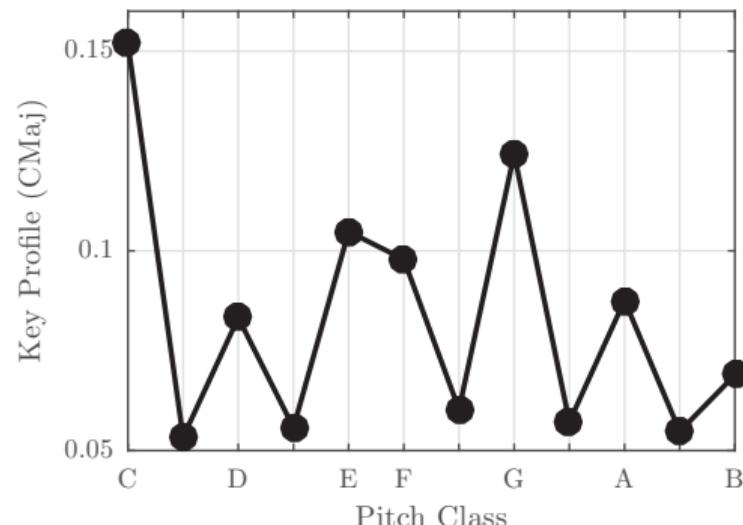
- ① define reference distribution for specific keys
- ② extract average pitch chroma from audio
- ③ compute distance between template and extracted chroma



# key detection

processing steps of simple key detection

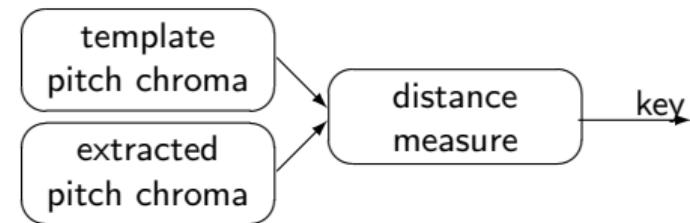
- ① define reference distribution for specific keys
- ② extract average pitch chroma from audio
- ③ compute distance between template and extracted chroma



# key detection

processing steps of simple key detection

- ① define reference distribution for specific keys
- ② extract average pitch chroma from audio
- ③ compute distance between template and extracted chroma



overview  
o

key  
oooo

pitch chroma  
oooooooo

key detection  
oo●oooo

summary  
o

# key detection

key template distance animation



matlab source: [matlab/animateKeyDetection.m](#)



# key detection

## key templates 1/2

- *Orthogonal  $\nu_o$* : root note is most salient component, other components negligible
  - same distance to all keys
  - no distinction between major and minor
- *Smoothed Orthogonal  $\nu_s$* : root note most salient, neighboring components less important
  - increasing key distance to tritone
  - no real distinction between major and minor
- *Diatonic  $\nu_d$* : all key-inherent pitches weighted equally
  - linear increasing key distance
- *Probe tone Ratings  $\nu_p$* : derived from perceptual tonal similarity
- *Extracted Key Profiles  $\nu_t$* : derived from real-world data

# key detection

## key templates 1/2

- *Orthogonal  $\nu_o$* : root note is most salient component, other components negligible
  - same distance to all keys
  - no distinction between major and minor
- *Smoothed Orthogonal  $\nu_s$* : root note most salient, neighboring components less important
  - increasing key distance to tritone
  - no real distinction between major and minor
- *Diatonic  $\nu_d$* : all key-inherent pitches weighted equally
  - linear increasing key distance
- *Probe tone Ratings  $\nu_p$* : derived from perceptual tonal similarity
- *Extracted Key Profiles  $\nu_t$* : derived from real-world data

# key detection

## key templates 1/2

- *Orthogonal  $\nu_o$* : root note is most salient component, other components negligible
  - same distance to all keys
  - no distinction between major and minor
- *Smoothed Orthogonal  $\nu_s$* : root note most salient, neighboring components less important
  - increasing key distance to tritone
  - no real distinction between major and minor
- *Diatonic  $\nu_d$* : all key-inherent pitches weighted equally
  - linear increasing key distance
- *Probe tone Ratings  $\nu_p$* : derived from perceptual tonal similarity
- *Extracted Key Profiles  $\nu_t$* : derived from real-world data

# key detection

## key templates 1/2

- *Orthogonal  $\nu_o$* : root note is most salient component, other components negligible
  - same distance to all keys
  - no distinction between major and minor
- *Smoothed Orthogonal  $\nu_s$* : root note most salient, neighboring components less important
  - increasing key distance to tritone
  - no real distinction between major and minor
- *Diatonic  $\nu_d$* : all key-inherent pitches weighted equally
  - linear increasing key distance
- *Probe tone Ratings  $\nu_p$* : derived from perceptual tonal similarity
- *Extracted Key Profiles  $\nu_t$* : derived from real-world data

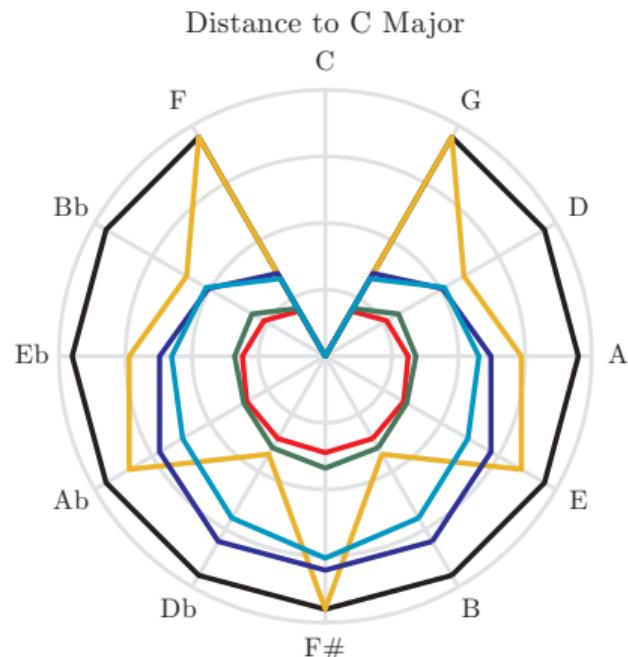
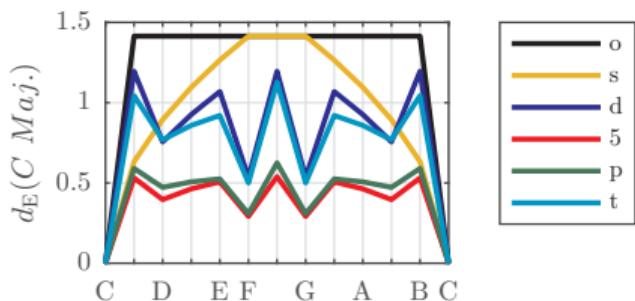
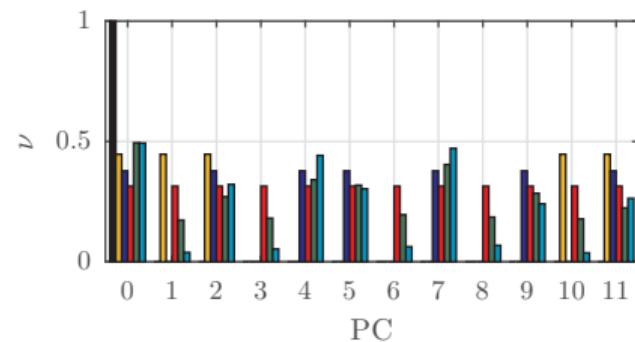
# key detection

## key templates 1/2

- *Orthogonal  $\nu_o$* : root note is most salient component, other components negligible
  - same distance to all keys
  - no distinction between major and minor
- *Smoothed Orthogonal  $\nu_s$* : root note most salient, neighboring components less important
  - increasing key distance to tritone
  - no real distinction between major and minor
- *Diatonic  $\nu_d$* : all key-inherent pitches weighted equally
  - linear increasing key distance
- *Probe tone Ratings  $\nu_p$* : derived from perceptual tonal similarity
- *Extracted Key Profiles  $\nu_t$* : derived from real-world data

# key detection

## key templates 2/2



# key detection

## variants

- **tonalness weight:**

estimate the tonality/noisiness and weight instantaneous pitch chroma

- **multiple estimations:**

split piece into regions and estimate key through majority

- **real-time key detection:**

estimate in sliding window

# key detection

## variants

- **tonalness weight:**  
estimate the tonality/noisiness and weight instantaneous pitch chroma
- **multiple estimations:**  
split piece into regions and estimate key through majority
- **real-time key detection:**  
estimate in sliding window

# key detection

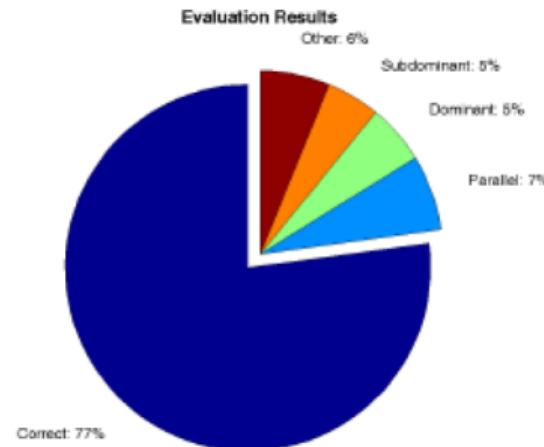
## variants

- **tonalness weight:**  
estimate the tonality/noisiness and weight instantaneous pitch chroma
- **multiple estimations:**  
split piece into regions and estimate key through majority
- **real-time key detection:**  
estimate in sliding window

# key detection

## results & typical errors

- typical errors: related keys
  - Dominant
  - Subdominant
  - Relative
  - Major/Minor



graph from<sup>1</sup>

<sup>1</sup> A. Lerch, "Ein Ansatz zur automatischen Erkennung der Tonart in Musikdateien," in *Proceedings of the VDT International Audio Convention (23. Tonmeistertagung)*, Leipzig, Nov. 2004.

# summary

## lecture content

### ● musical key

- set of pitch classes constructing pitched content
- defined by *tonic* (important center) and *mode* (scale)

### ● pitch chroma

- reduced 12-dimensional octave-independent pitch representation
- relatively robust against timbre variation

### ● automatic key recognition

- standard approach is template-based
- extracted average pitch chroma is compared with predefined template
- inverse distance measure indicates key likelihoods

