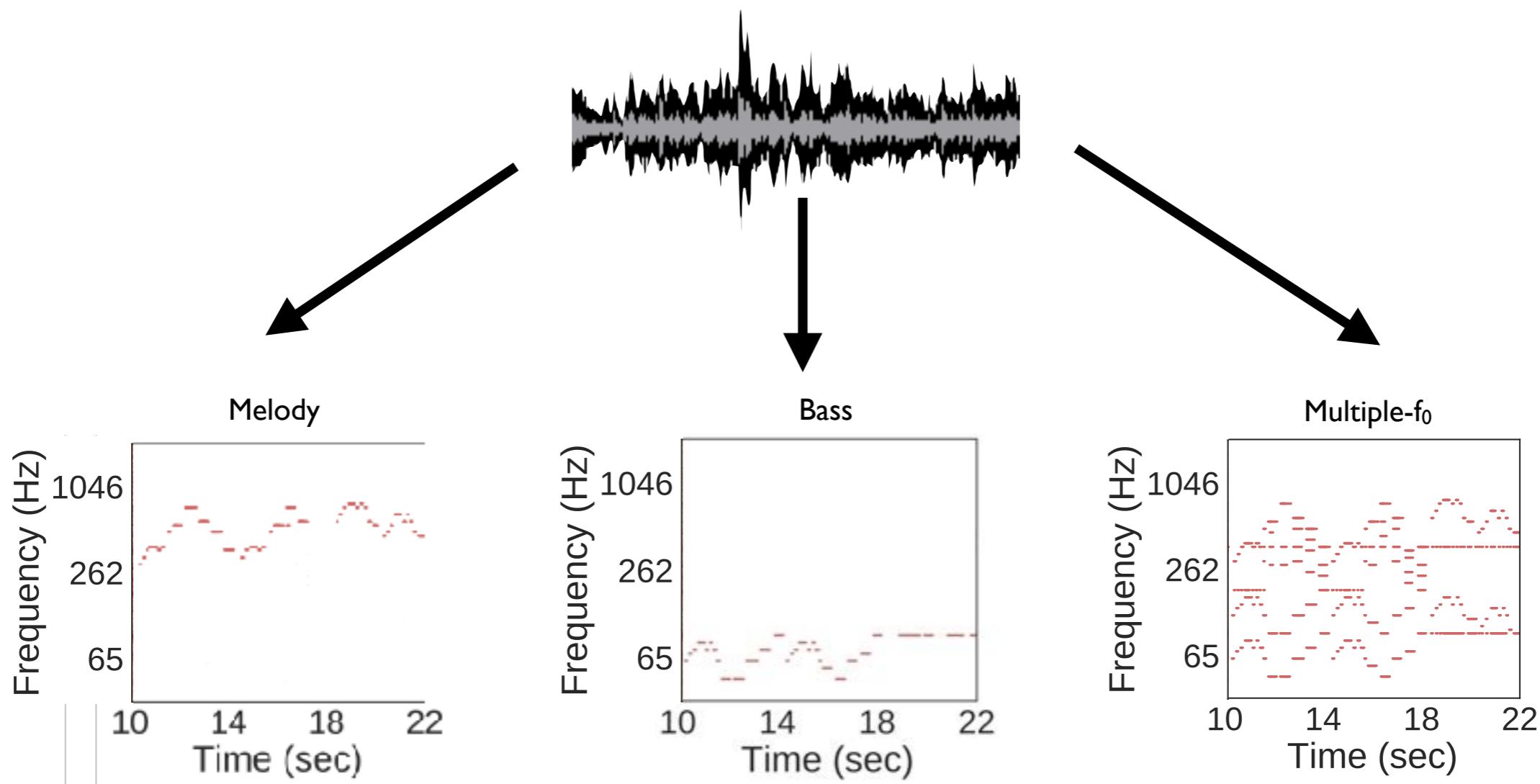


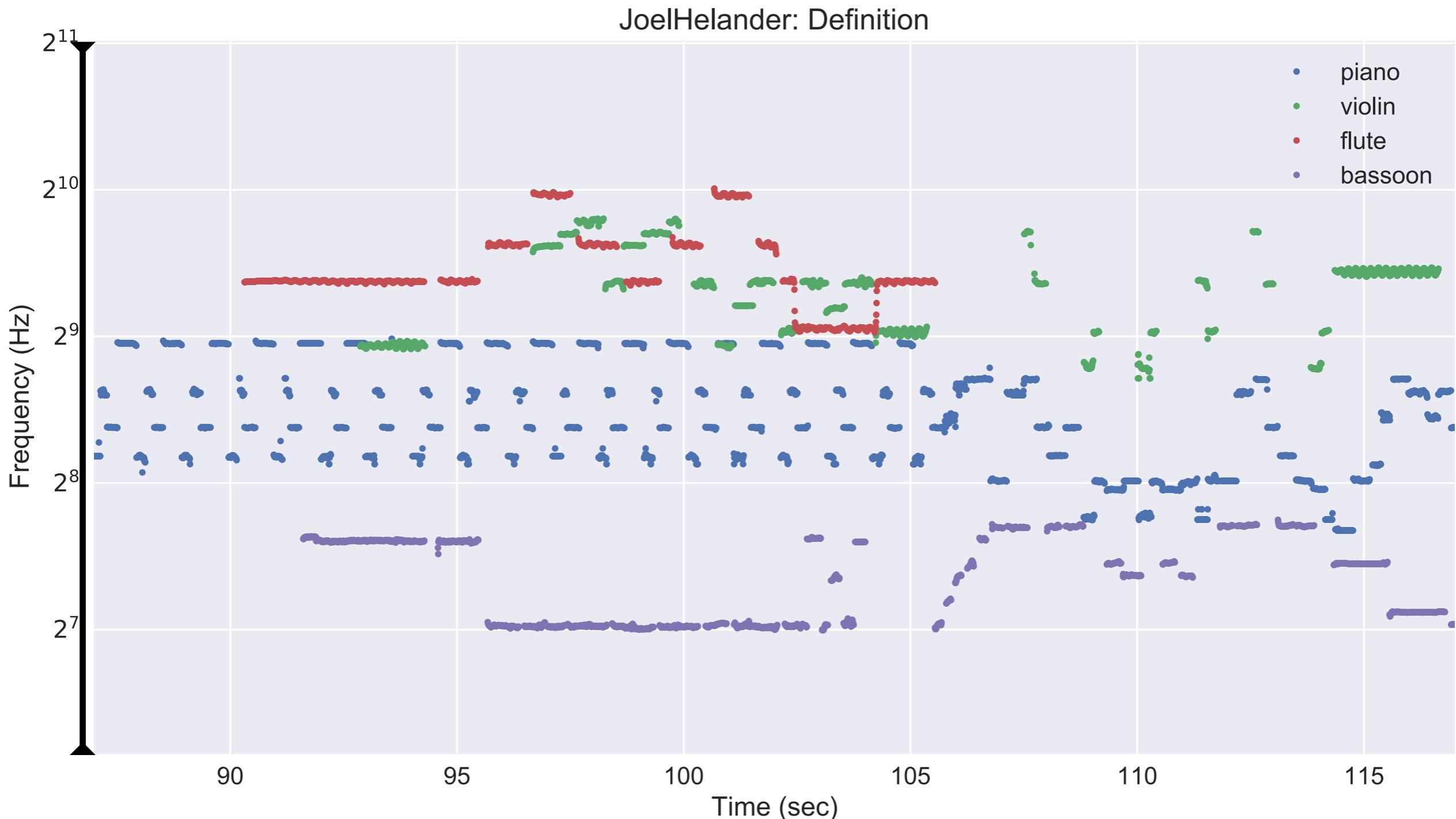
# Polyphonic Fundamental Frequency Estimation

$f_0$  Estimation Tutorial - Part 2

# Polyphonic $f_0$ Estimation



# Polyphonic $f_0$ Estimation



# Polyphonic $f_0$ estimation tasks

- Multiple- $f_0$ /multi-pitch estimation



- Melody estimation



- Vocal- $f_0$  estimation (not the same!)



- Bass line estimation



- (solo) polyphonic source transcription



- source-specific  $f_0$  estimation



# Applications

- Transcription



- Music Generation + Style Transfer



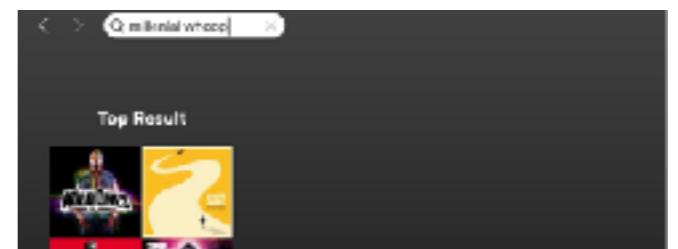
- Audio Effects



- Performance Analysis



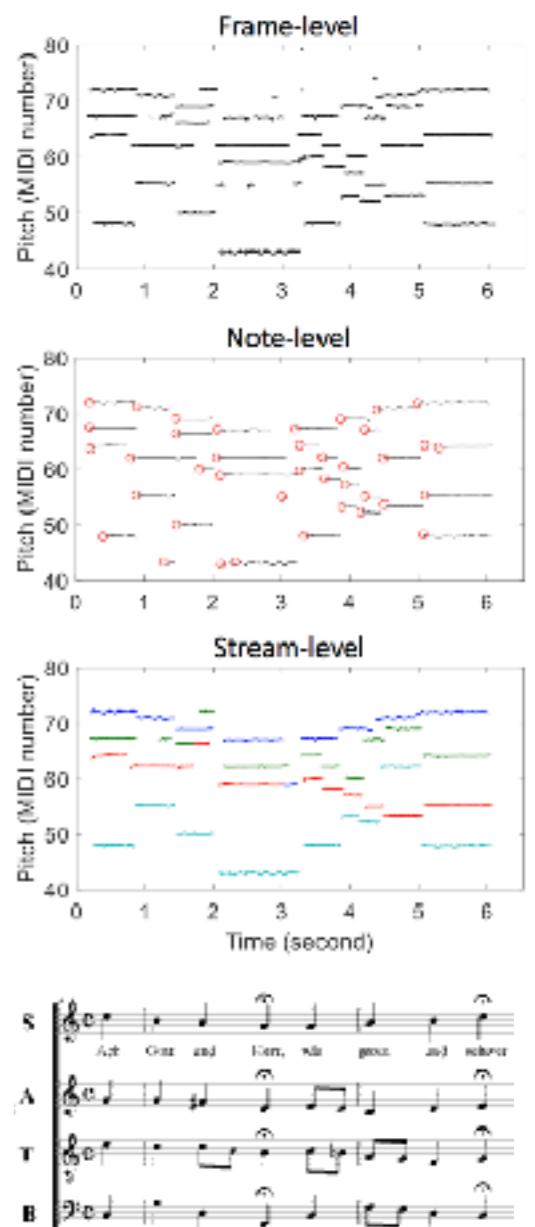
- Similarity and Retrieval



# Background

# Related Tasks

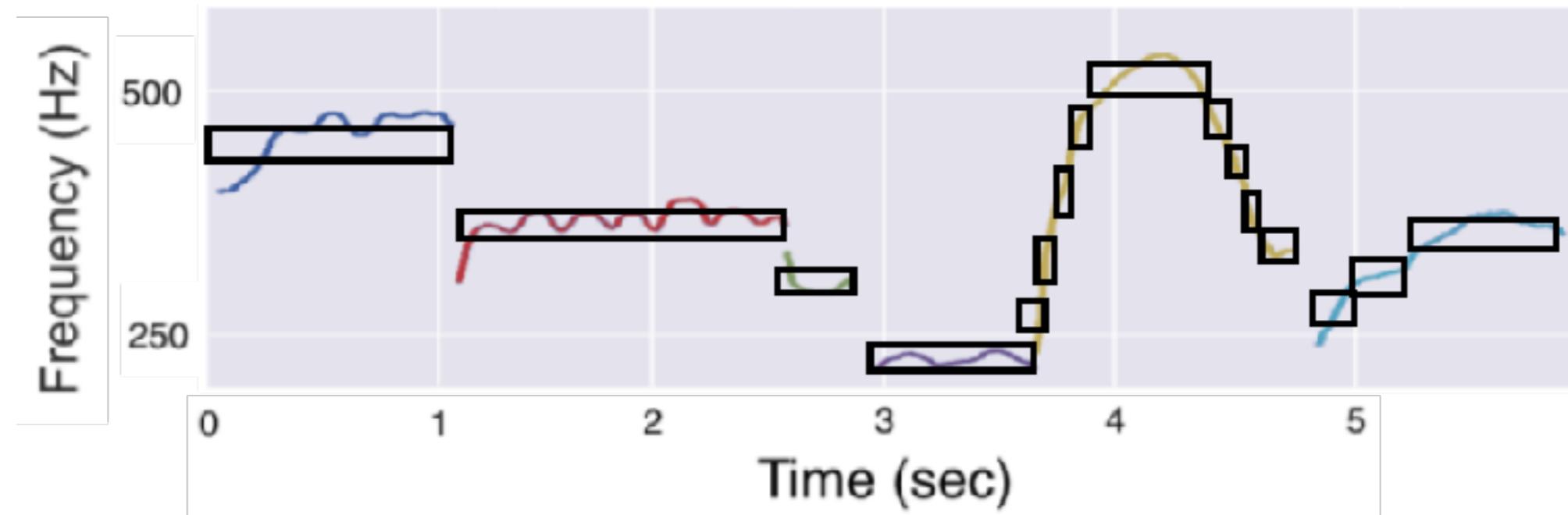
- frame-level  $f_0$  estimation
- Note estimation
- Streaming
- Score transcription



Emmanouil Benetos, Simon Dixon, Zhiyao Duan, and Sebastian Ewert  
“Automatic Music Transcription: An Overview”  
To appear – IEEE Signal Processing Magazine

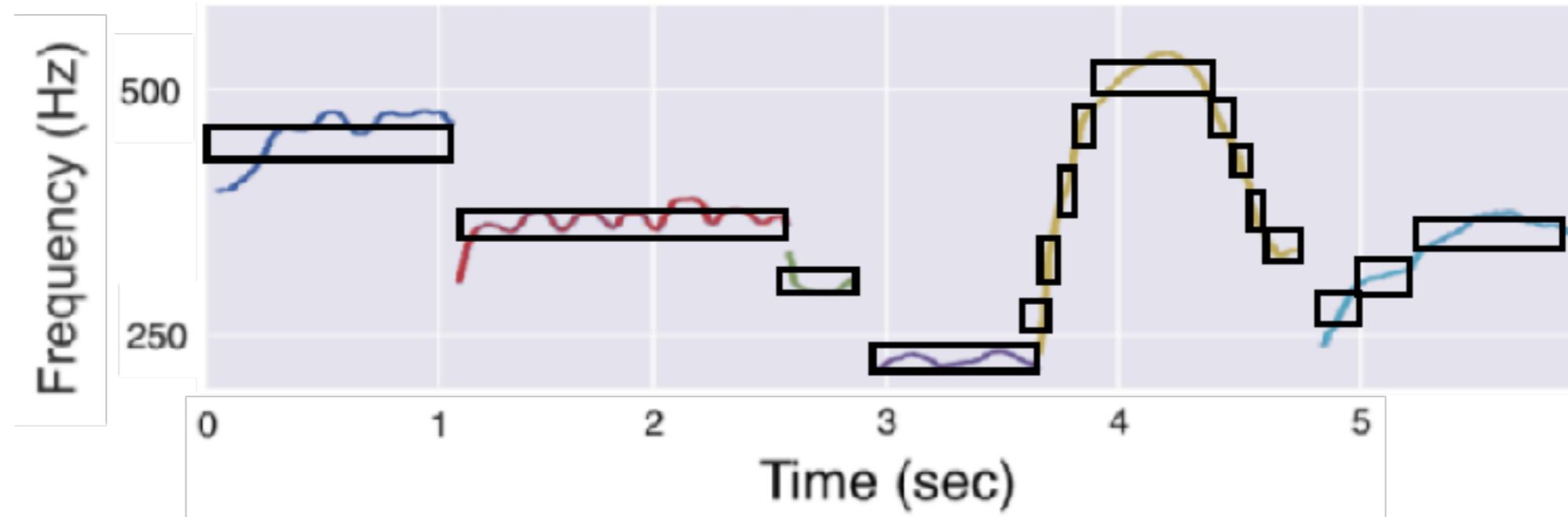


# Time Resolution



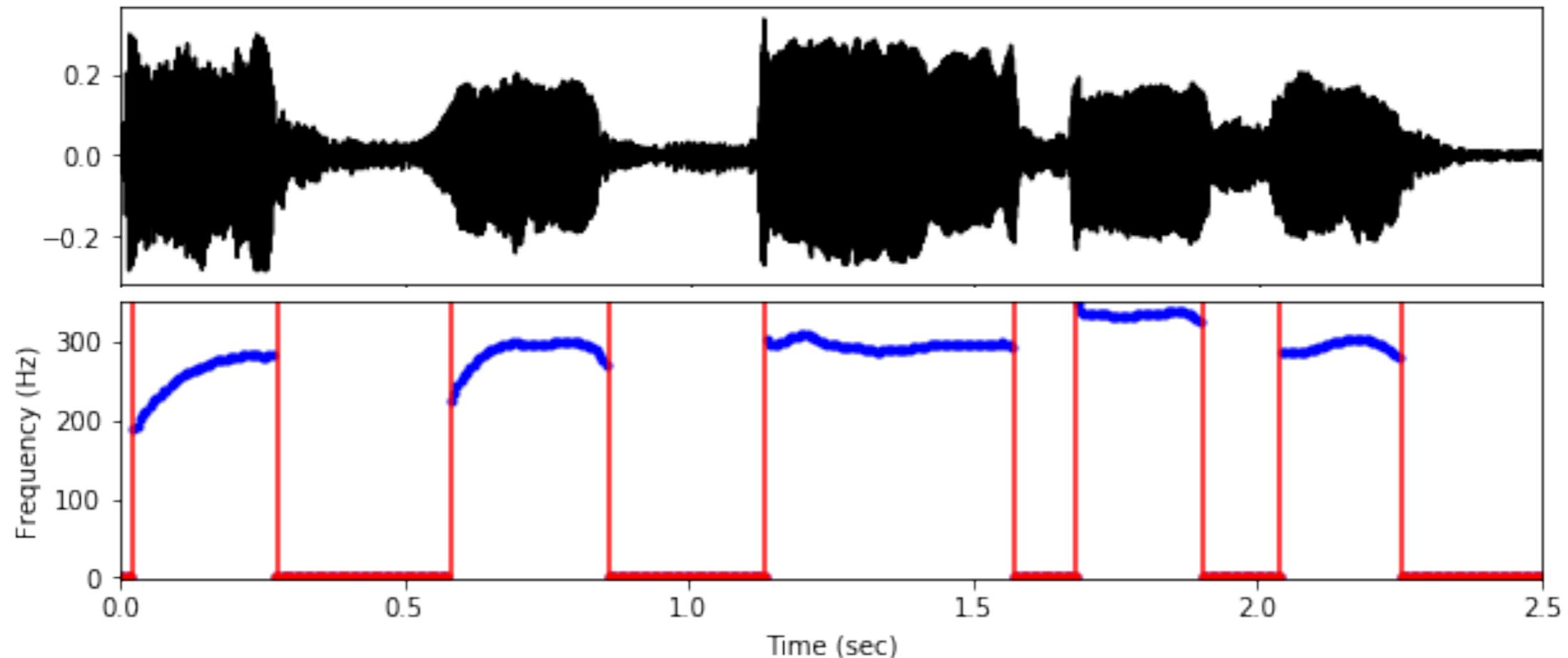
- **Fine:**  $f_0$  trajectories - (time,  $f_0$ )
- **Coarse:** Notes - (start time, end time,  $f_0$ )
- **Coarsest:** Quantized Notes - times are quantized to fractional beats

# Frequency Resolution



- **Discrete:** Semitone resolution (integer midi notes)
- **“Continuous”:** less than semitone resolution (e.g. 1/5th semitone)

# Voicing

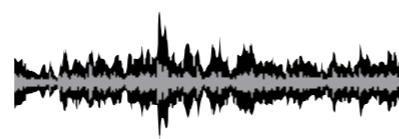


- “**Voiced**” - represented with  $f_0$  value  $> 0$
- “**Unvoiced**” - represented with values  $\leq 0$

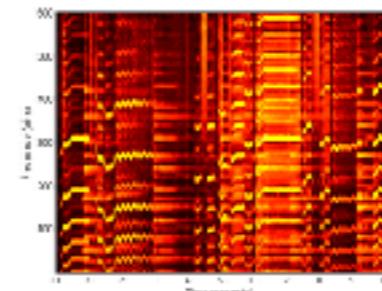
# Approaches to $f_0$ Estimation

# Typical Approach

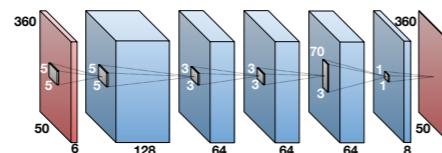
Polyphonic Music



Time Frequency Transform

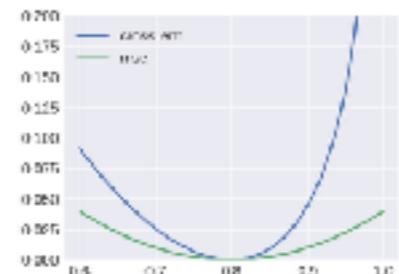


Model



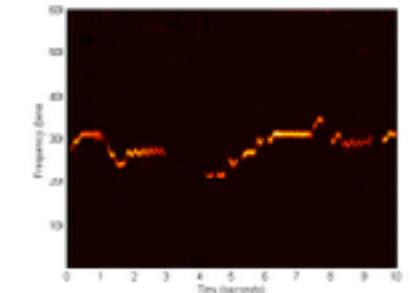
+

Optimization Function

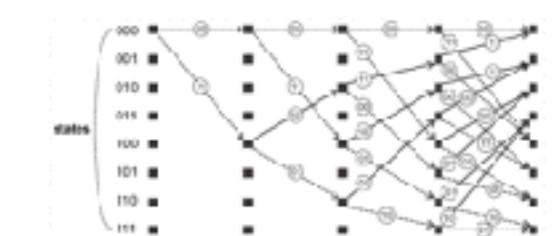
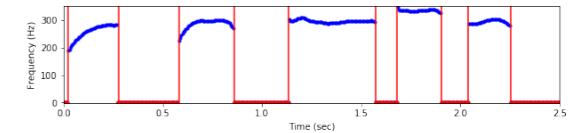


Polyphonic  $f_0$  Estimation

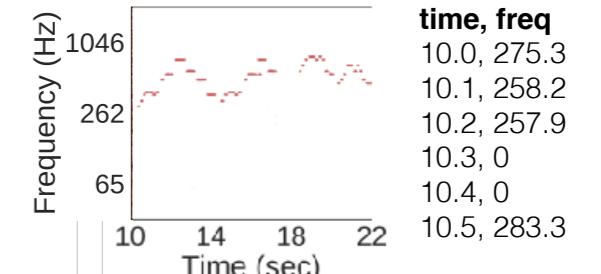
Salience Representation



Voicing Determination + Decoding



$f_0$  Time Series

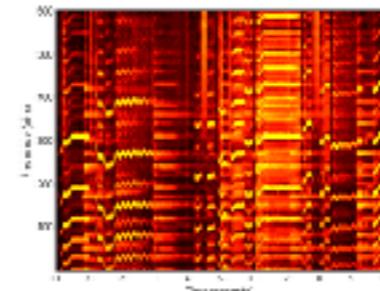


# Typical Approach

Polyphonic Music



Time Frequency Transform

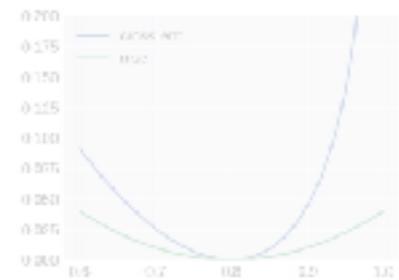


Model



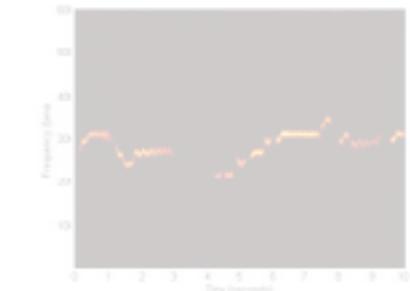
+

Optimization Function

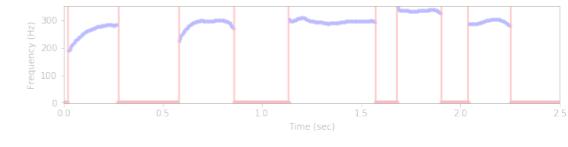


Polyphonic  $f_0$  Estimation

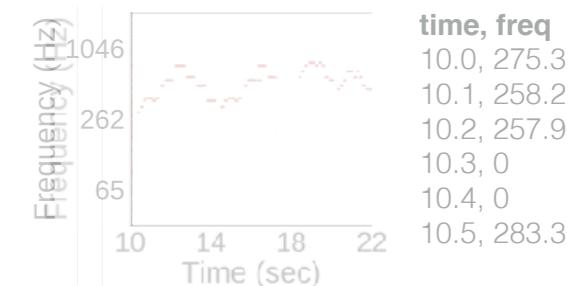
Salience Representation



Voicing Determination + Decoding



$f_0$  Time Series



# Time-Frequency Transforms

Magnitude STFT

Matija Marolt

"On Finding Melodic Lines in Audio Recordings"  
DAFx (2004)

Multi-resolution FFT

Karin Dressler

"Towards Computational Auditory Scene Analysis: Melody Extraction from Polyphonic Music"  
CMMR (2012)

Magnitude CQT

Benoit Fuentes, Antoine Liutkus, Roland Badeau, Gaël Richard

"Probabilistic Model for Main Melody Extraction using constant-Q transform"  
ICASSP (2012)

HCQT

Rachel Bittner, Brian McFee, Justin Salamon, Peter Li, Juan Pablo Bello

"Deep Salience Representations for F0 Estimation in Polyphonic Music"  
ISMIR (2017)

Justin Salamon, Emilia Gómez, Jordi Bonada

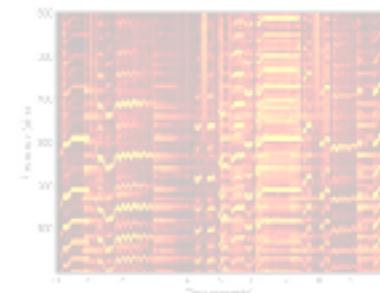
"Sinusoid Extraction and Salience Function Design for Predominant Melody Estimation"  
DAFx (2011)

# Typical Approach

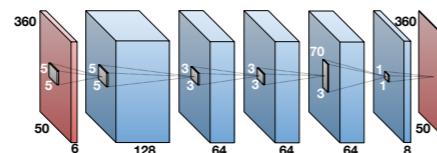
Polyphonic Music



Time Frequency Transform

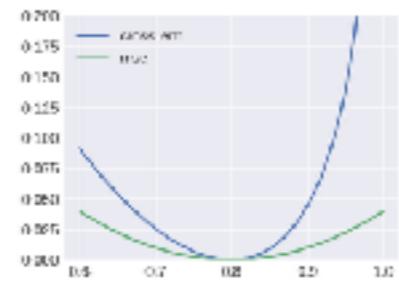


Model



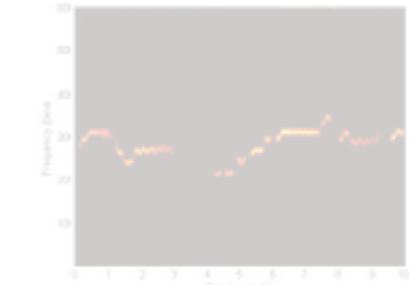
+

Optimization Function

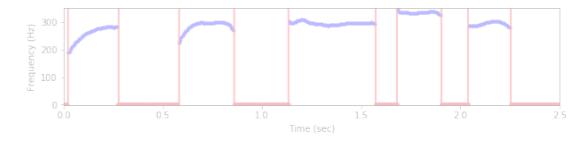


Polyphonic  $f_0$  Estimation

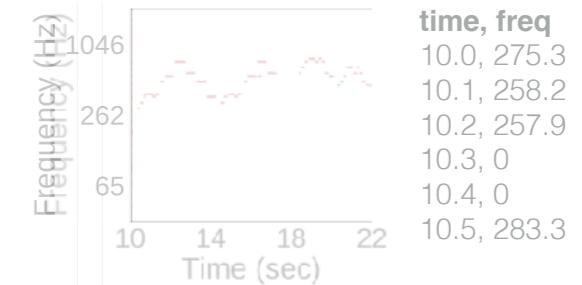
Salience Representation



Voicing Determination + Decoding



$f_0$  Time Series



# Models + Optimization

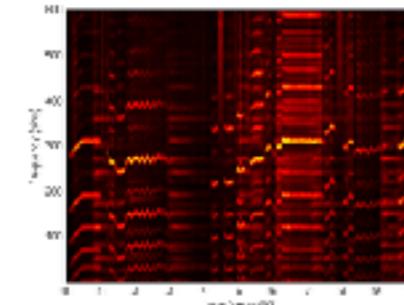
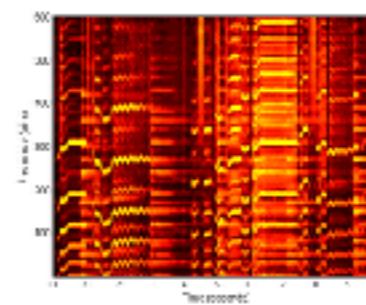
1) Salience Modeling

2) Nonnegative Matrix Factorization

3) Neural Networks

# Models + Optimization

## 1) Salience Modeling



## Harmonic Summation

Matti Ryynänen and Anssi Klapuri  
“Automatic transcription of melody, baseline and chords in polyphonic music”  
CMJ (2008)

## Source Filter Modeling

Jean Louis Durrieu, Gaël Richard, Bertrand David, Cédric Févotte  
“Source/filter model for unsupervised main melody extraction from polyphonic audio signals.”  
IEEE TASLP (2010)

# Models + Optimization

## 2) Nonnegative Matrix Factorization

Input      Templates      Activations      Reconstruction



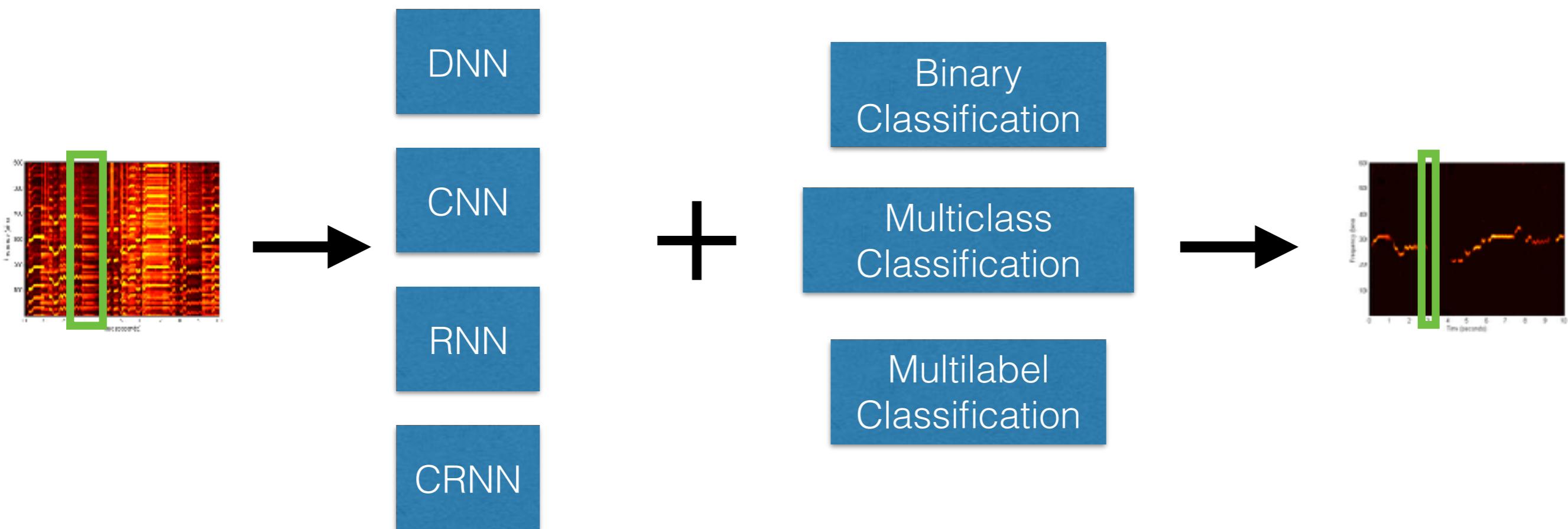
find the that minimize  $d(\text{Input}, \text{Reconstruction})$

Activations

Emmanouil Benetos, Simon Dixon, Zhiyao Duan, and Sebastian Ewert  
"Automatic Music Transcription: An Overview"  
To appear - IEEE Signal Processing Magazine

# Models + Optimization

## 3) Neural Networks

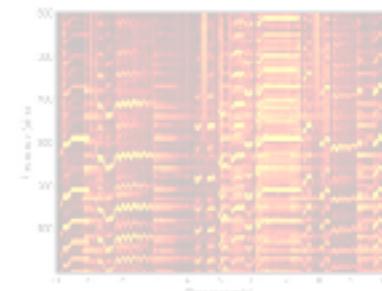


# Typical Approach

Polyphonic Music



Time Frequency Transform

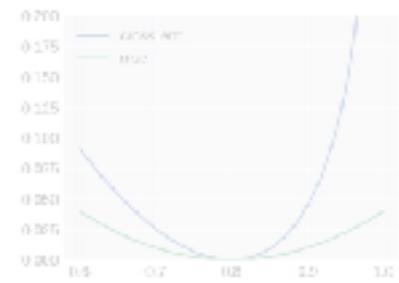


Model



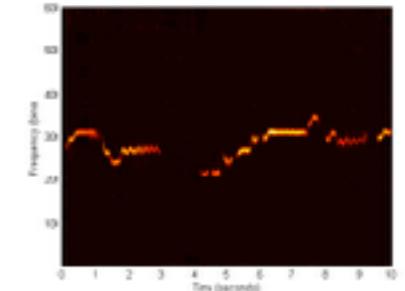
+

Optimization Function

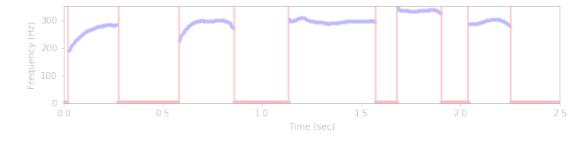


Polyphonic  $f_0$  Estimation

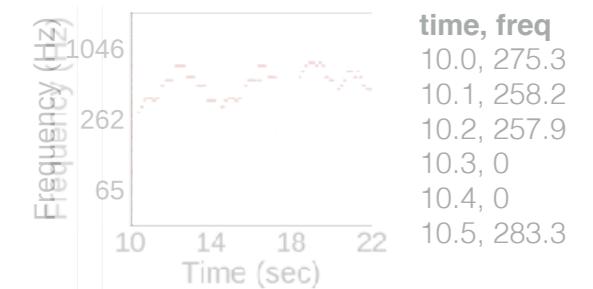
Salience Representation



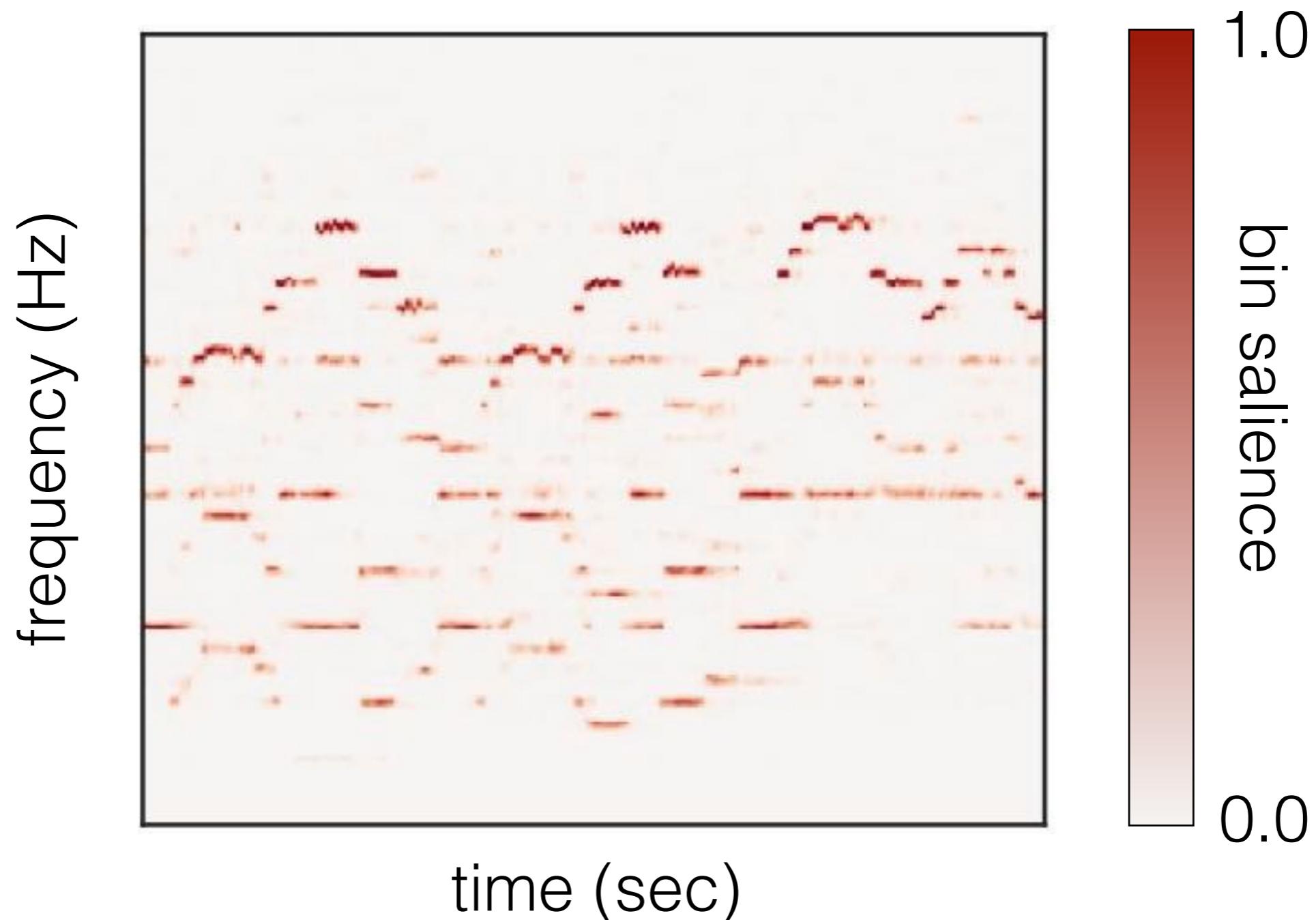
Voicing Determination + Decoding



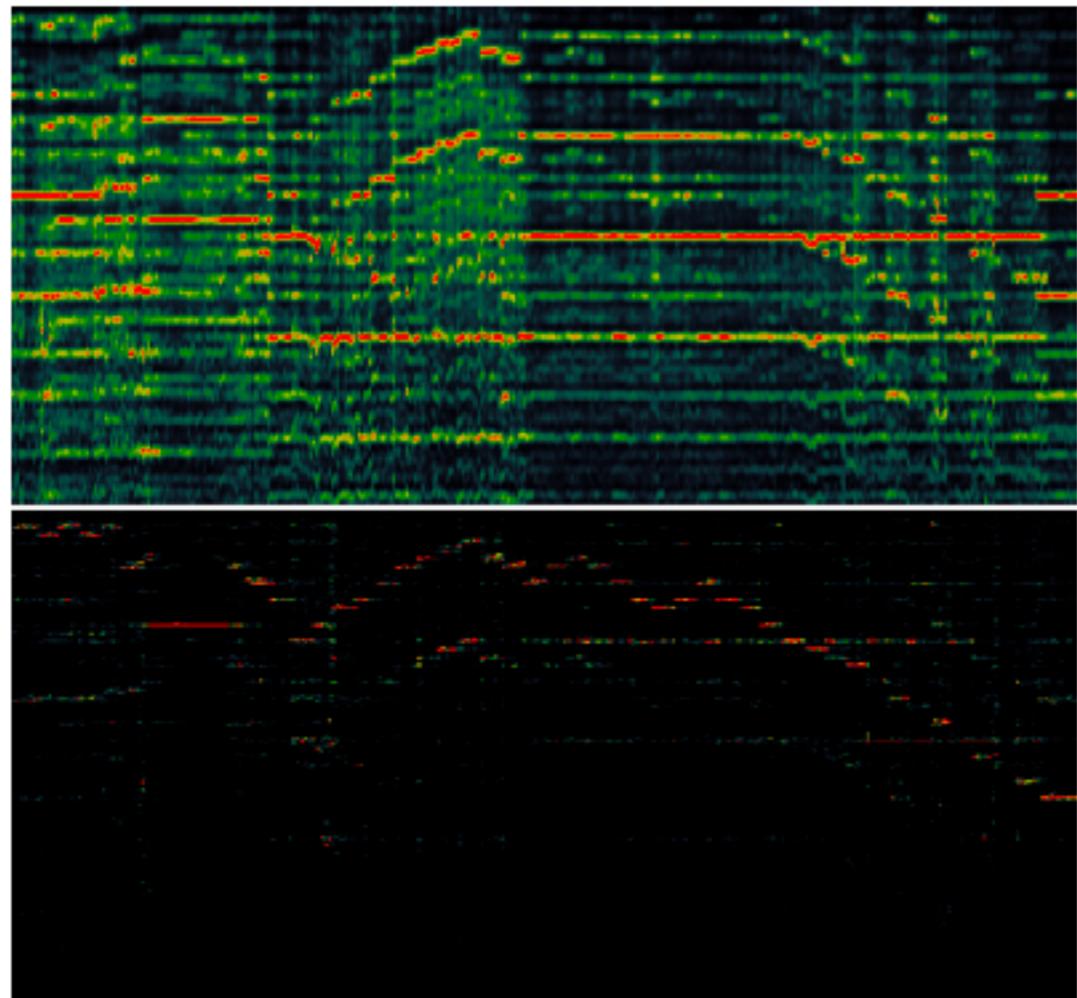
$f_0$  Time Series



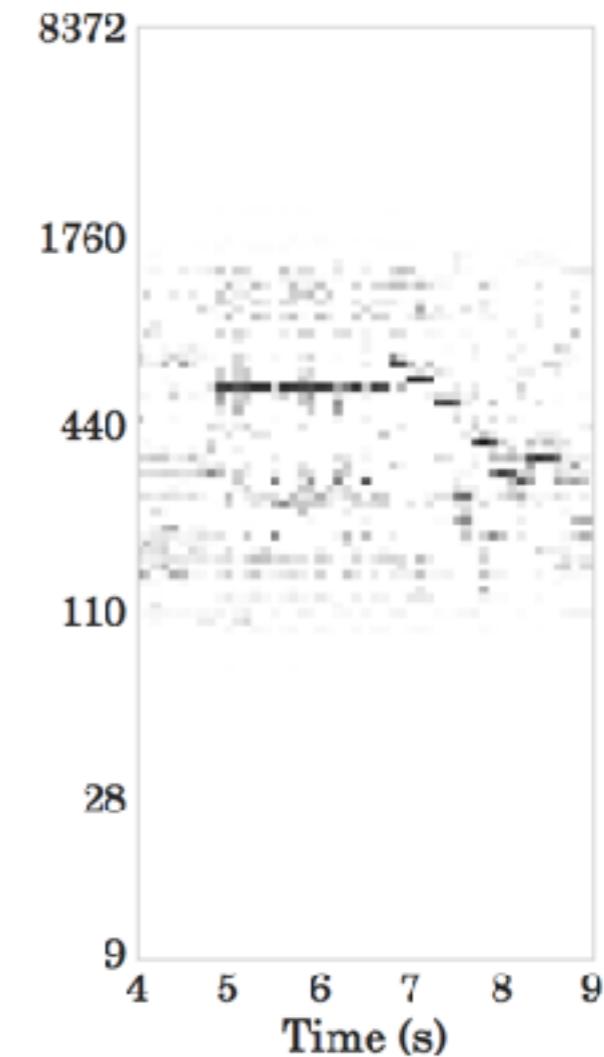
# Saliency Representations



# Saliency Representations



Juan José Bosch, Ricard Marxer, Emilia Gómez  
“Evaluation and Combination of Pitch  
Estimation Methods for Melody Extraction in  
Symphonic Classical Music”  
JNMR (2016)



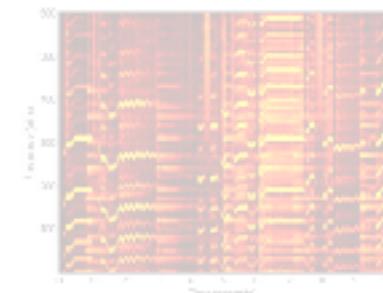
Stefan Balke, Christian Dittmar, Jakob  
Abeßer, Meinard Müller  
“Data Driven Solo Voice Enhancement for Jazz  
Music Retrieval”  
ICASSP (2017)

# Typical Approach

Polyphonic Music



Time Frequency Transform

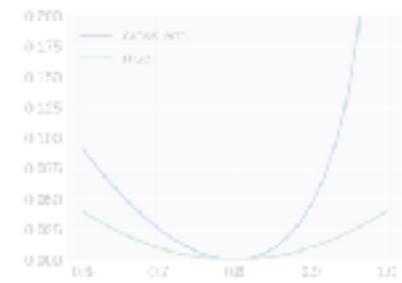


Model



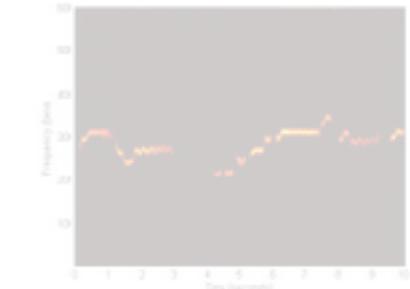
+

Optimization Function

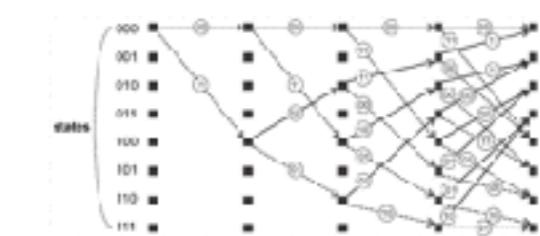
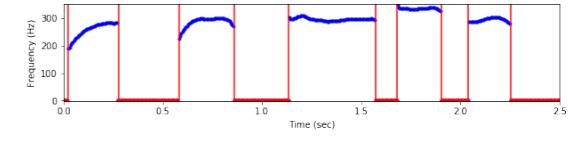


Polyphonic  $f_0$  Estimation

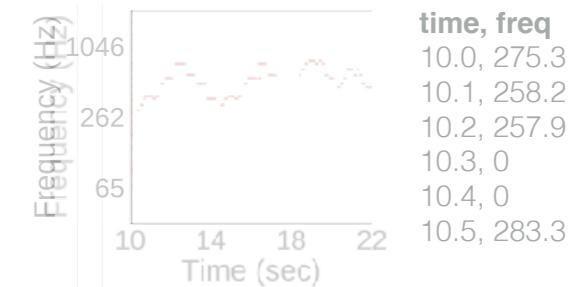
Salience Representation



Voicing Determination + Decoding



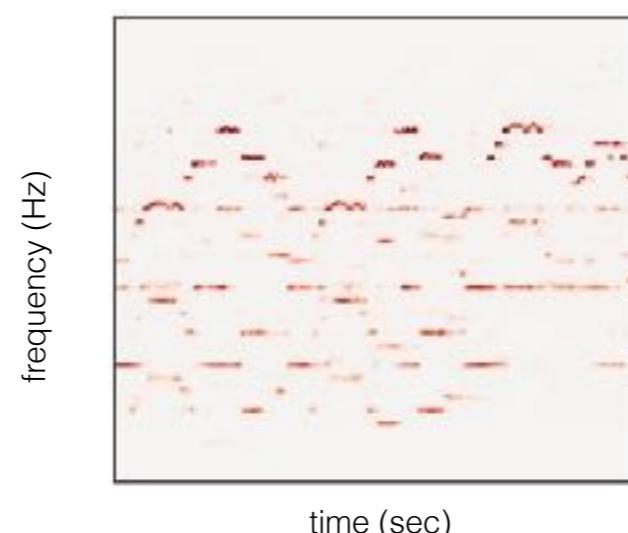
$f_0$  Time Series



# Voicing and Decoding

Thresholding

e.g. Vocal Activity  
Detection



argmax

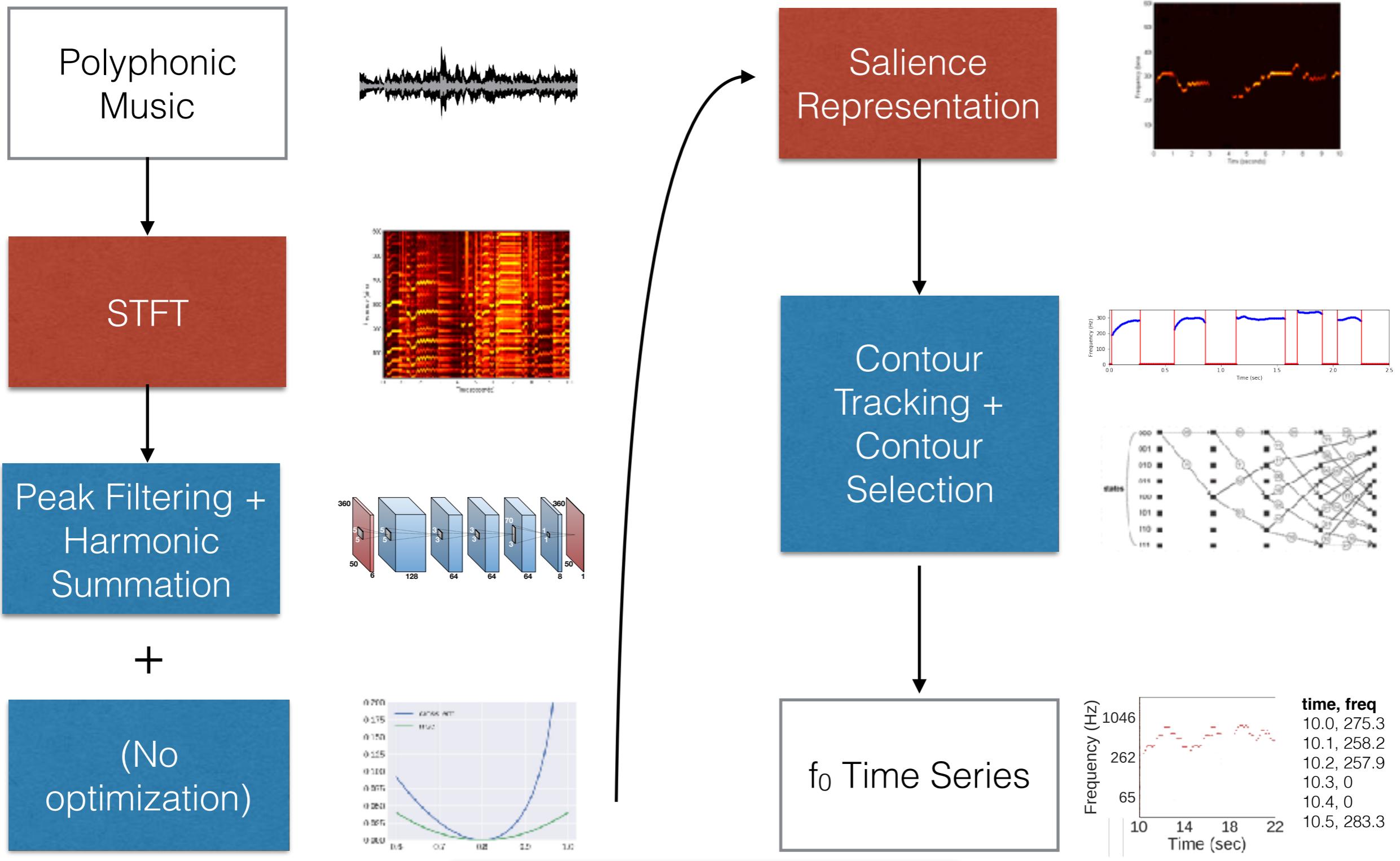
Viterbi

Contour Tracking

# Example 1: Salience

Justin Salamon, Emilia Gómez

"Melody Extraction from Polyphonic Music Signals using Pitch Contour Characteristics"  
IEEE TASLP (2012)

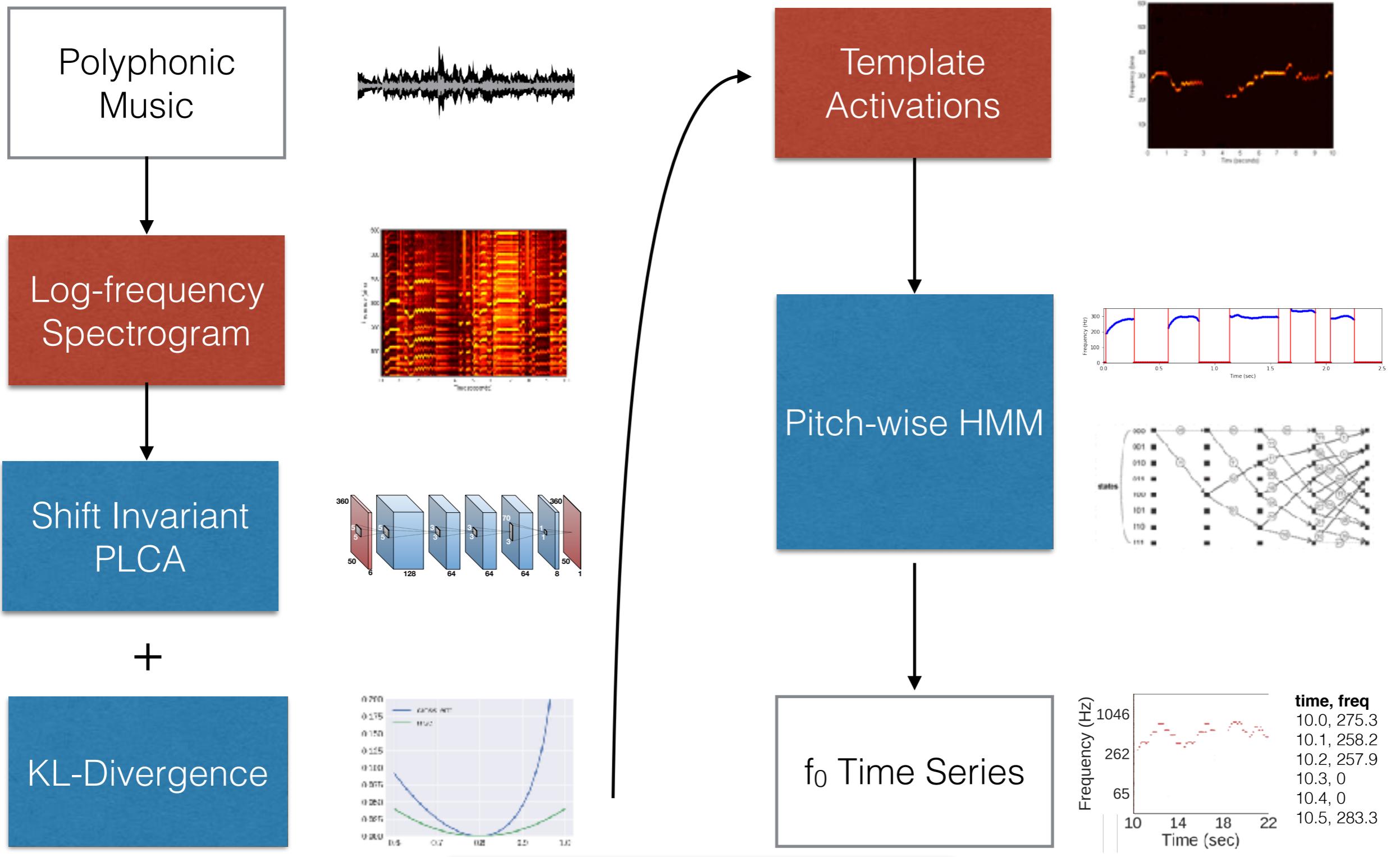


# Example 2: NMF

Emannouil Benetos, Simon Dixon

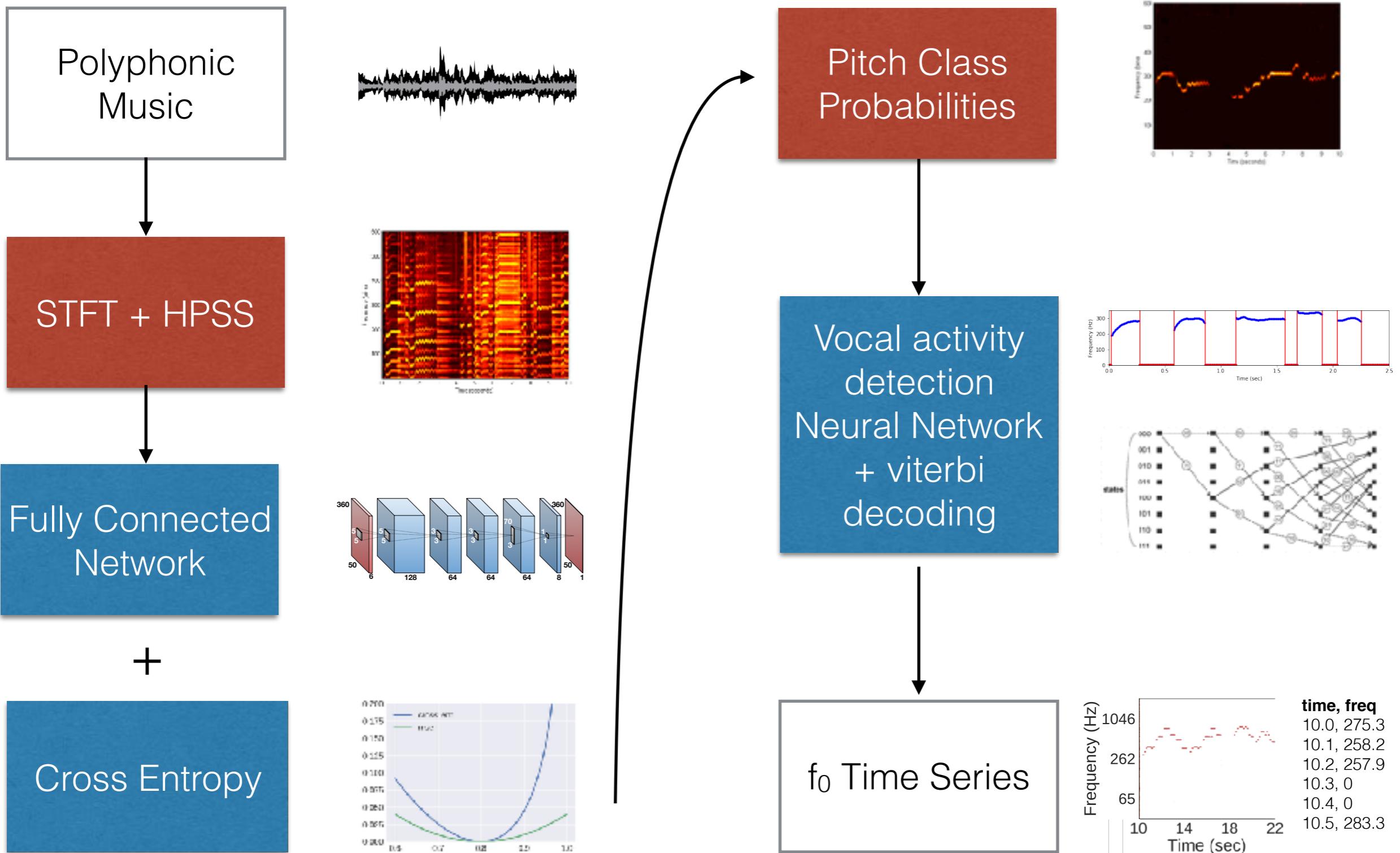
"Multiple-instrument polyphonic music transcription using a temporally constrained shift-invariant model"

JASA (2013)



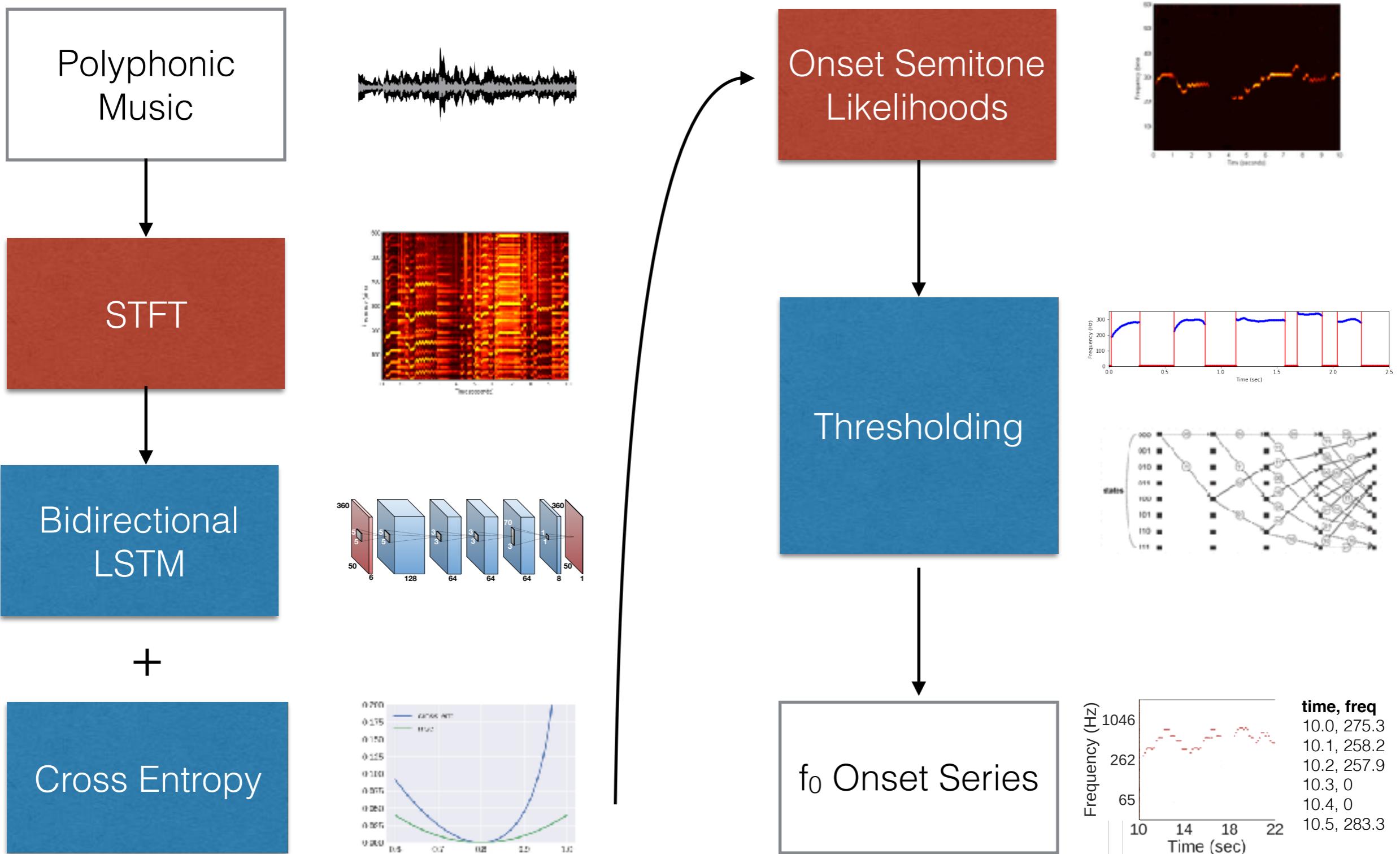
# Example 3: Neural Network 1

François Rigaud, Mattieu Radenac  
"Singing Voice Melody Transcription using Deep Neural Networks"  
ISMIR (2016)



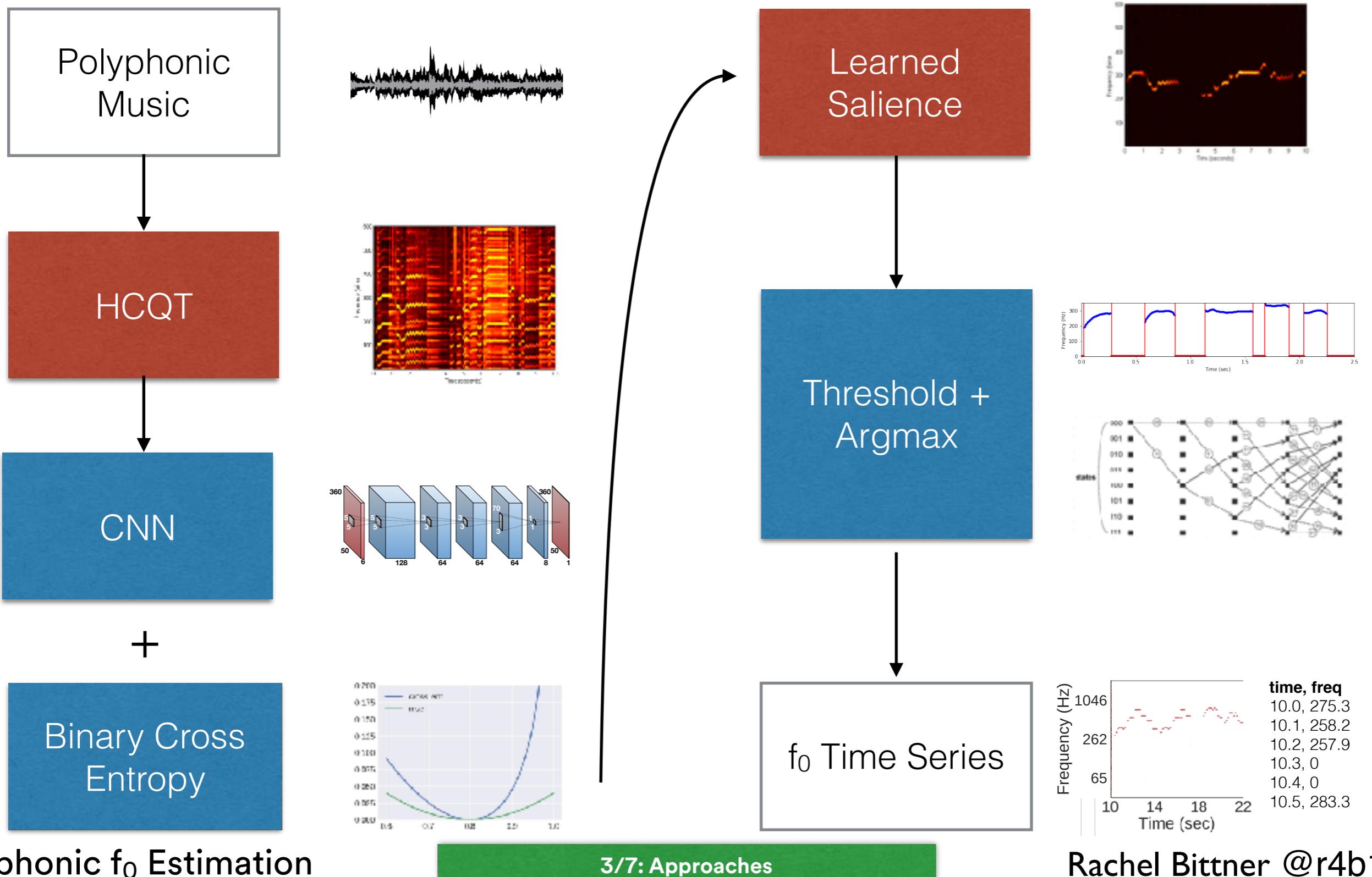
# Example 4: Neural Network 2

Sebastian Böck, Markus Schedl  
"Polyphonic piano note transcription with recurrent neural networks"  
ICASSP (2012)



# Example 5: Neural Network 3

Rachel Bittner, Brian McFee, Justin Salamon, Peter Li, Juan Pablo Bello  
"Deep Salience Representations for F0 Estimation in Polyphonic Music"  
ISMIR (2017)



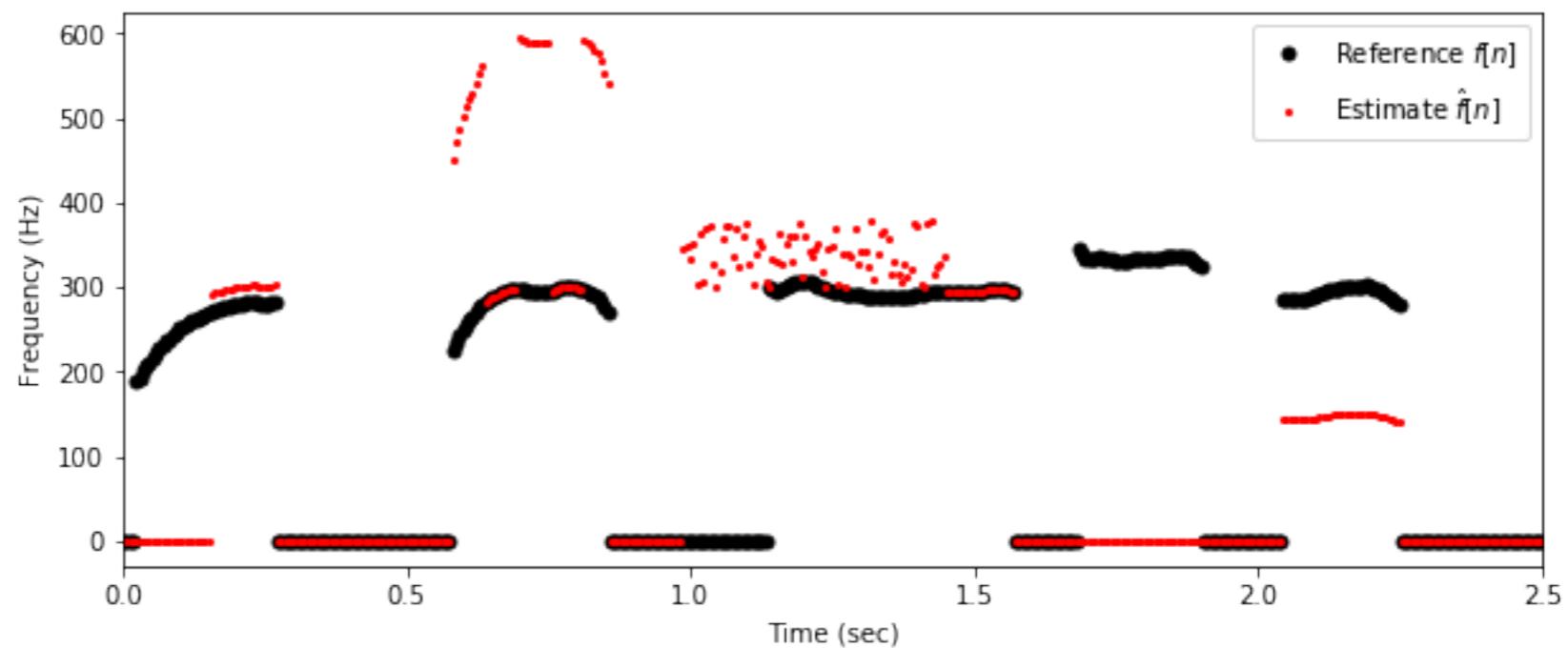
# What about source separation?



# Evaluation

# Single-f0 Qualitative Evaluation

## 1. Plots



## 2. Sonification

$$y(t) = a(t) \times \sin(2\pi \int^t f(x) dx)$$

>  $a(t)$  - amplitude over time

voicing: binary

confidence: continuous

>  $f(t)$  - frequency estimate



Mix Audio



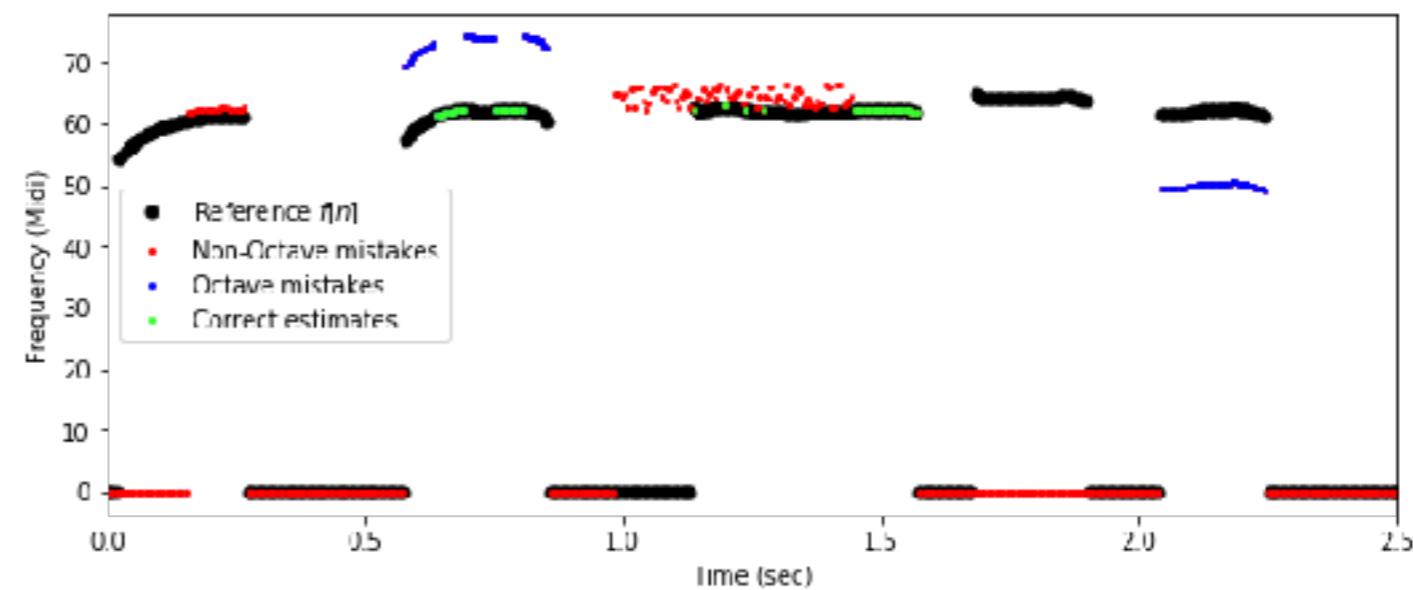
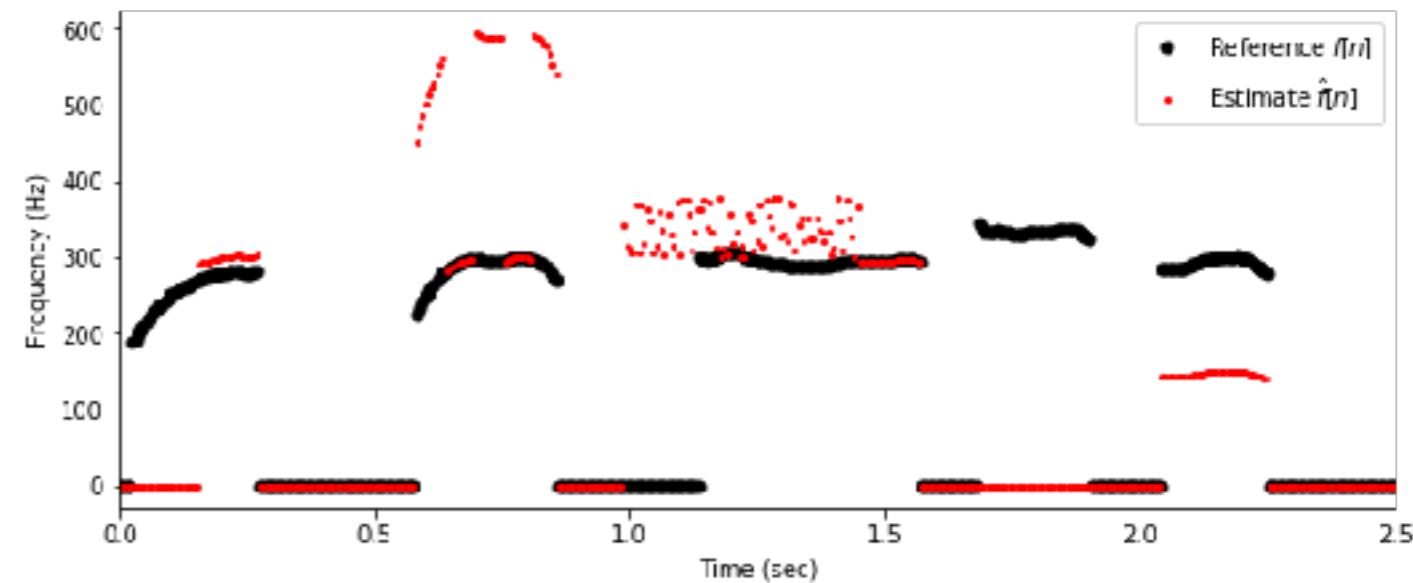
Reference f0



Estimate f0

# Single-f0 Quantitative Evaluation

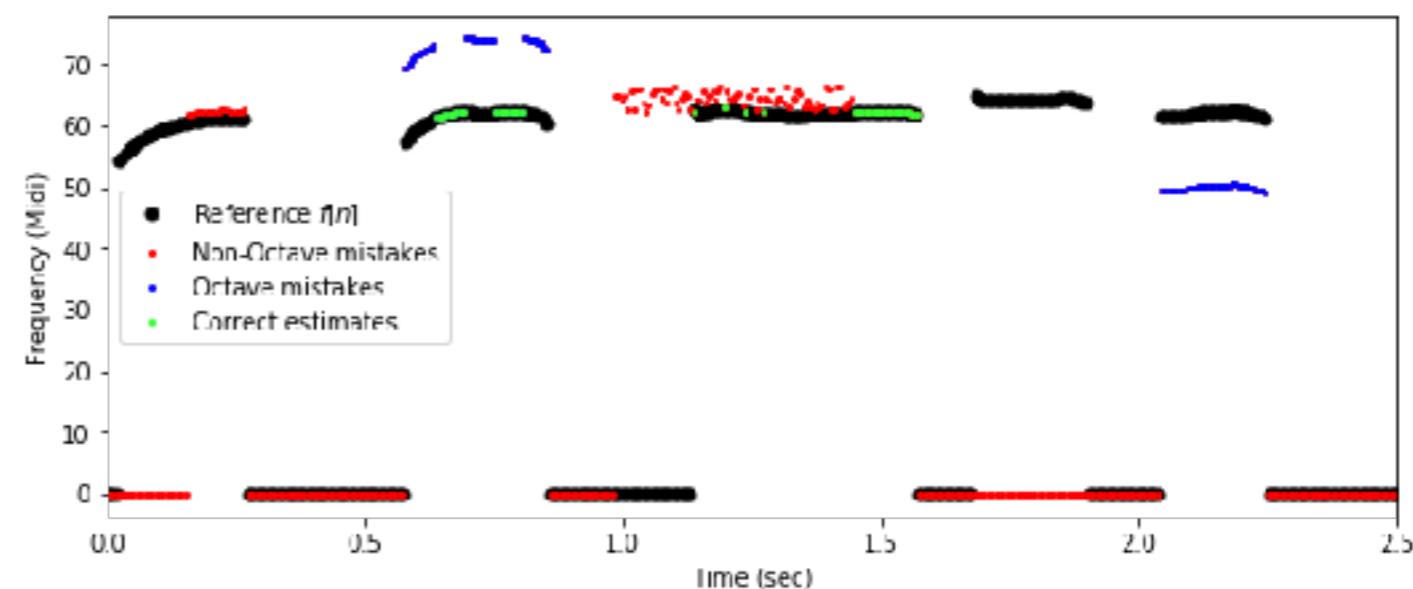
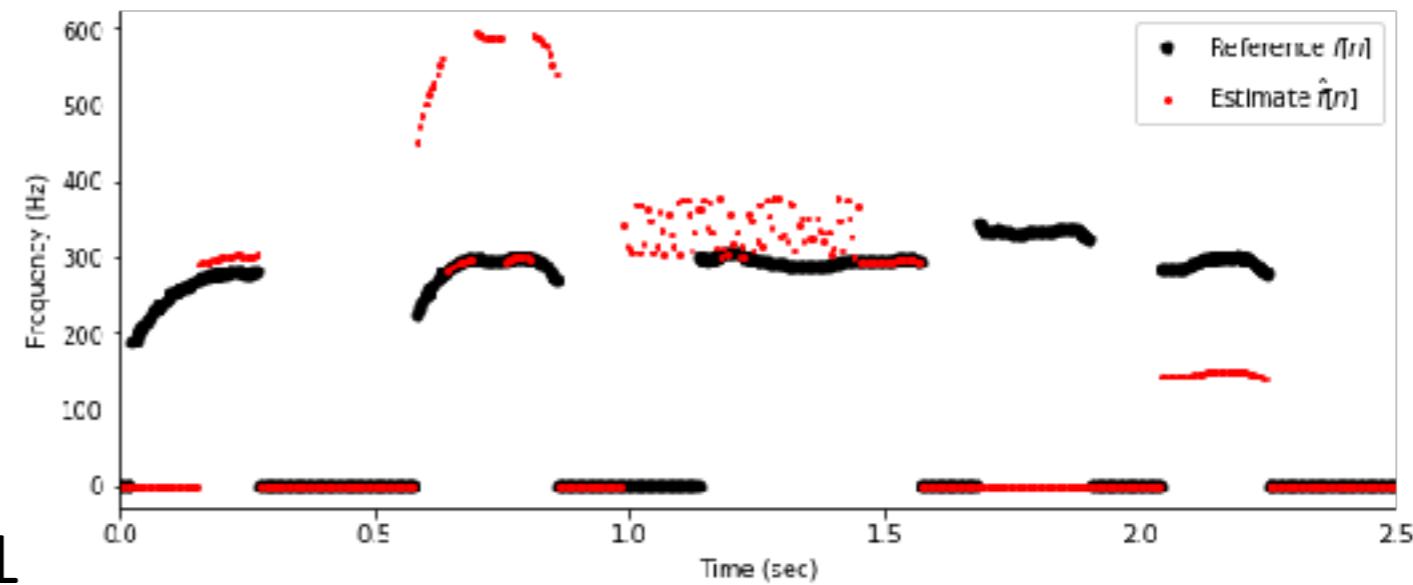
- Gross Percentage Error
- Mean Squared Error



# Single-f0 Quantitative Evaluation

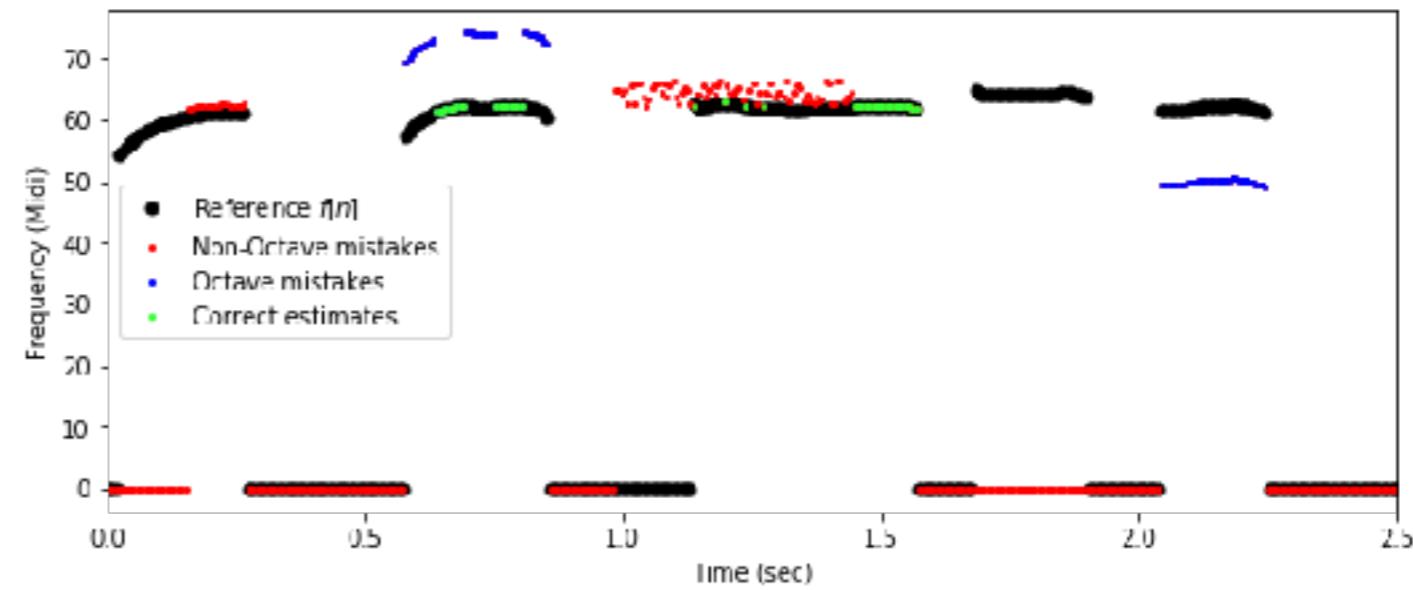
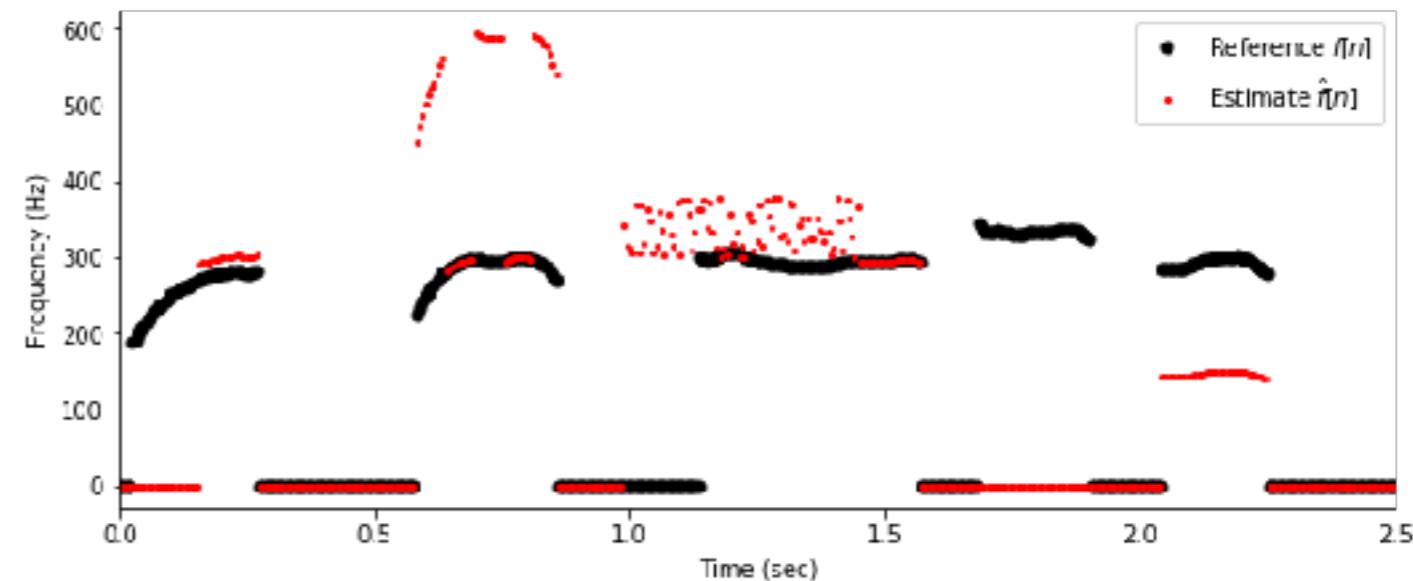
```
('Voicing Recall', 0.746)
('Voicing False Alarm', 0.136)
('Raw Pitch Accuracy', 0.1875)
('Raw Chroma Accuracy', 0.463)
('Overall Accuracy', 0.487)
```

[github.com/craffel/mir\\_eval](https://github.com/craffel/mir_eval)  
mir\_eval.melody.evaluate

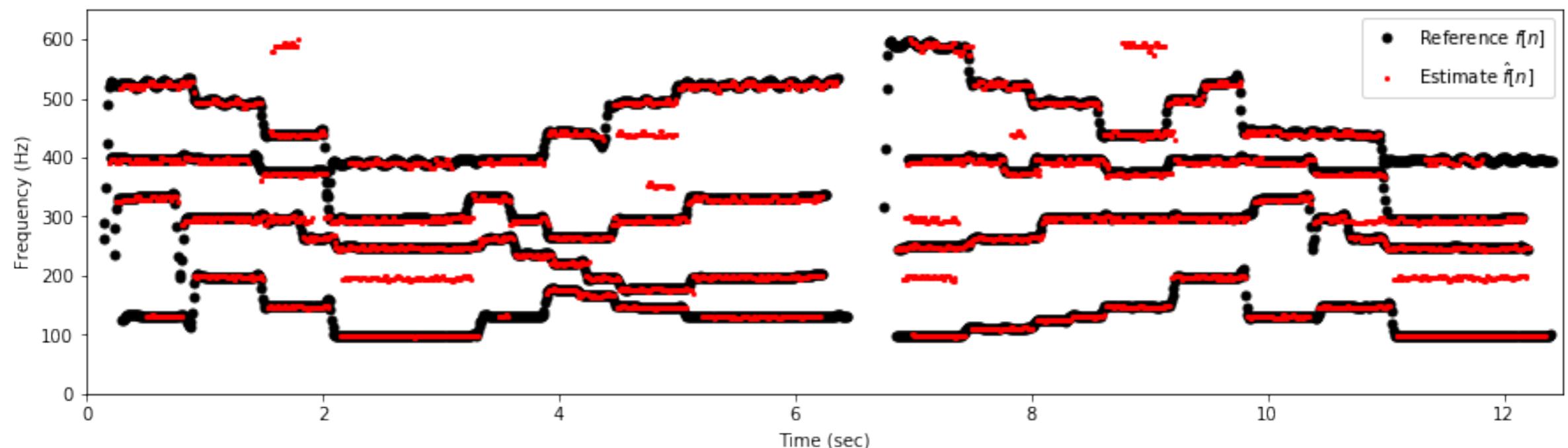


# Single-f0 Quantitative Evaluation

- Octave Jump Ratio (OJ)
- Chroma Continuity (CC)



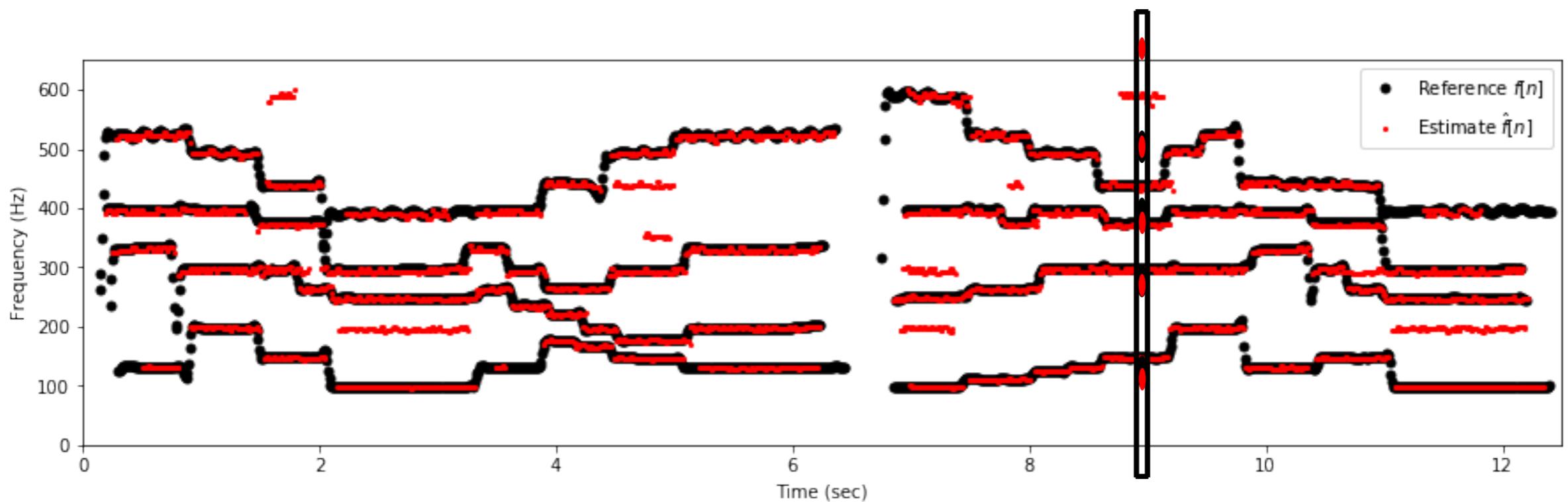
# Multiple-f0 Qualitative Evaluation



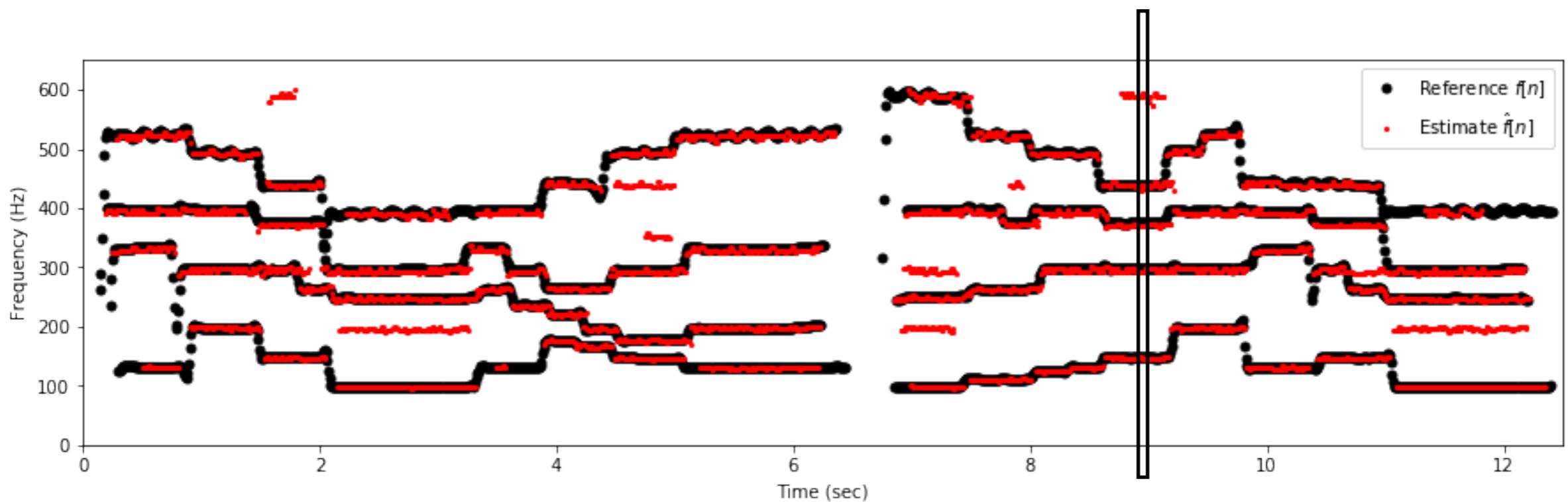
Reference f0

Estimate f0

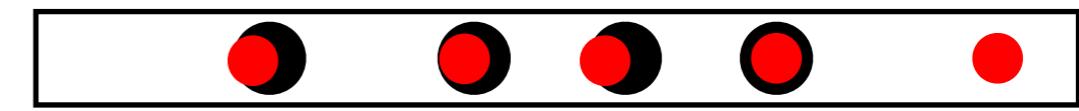
# Multiple-f0 Quantitative Evaluation



# Multiple-f0 Quantitative Evaluation



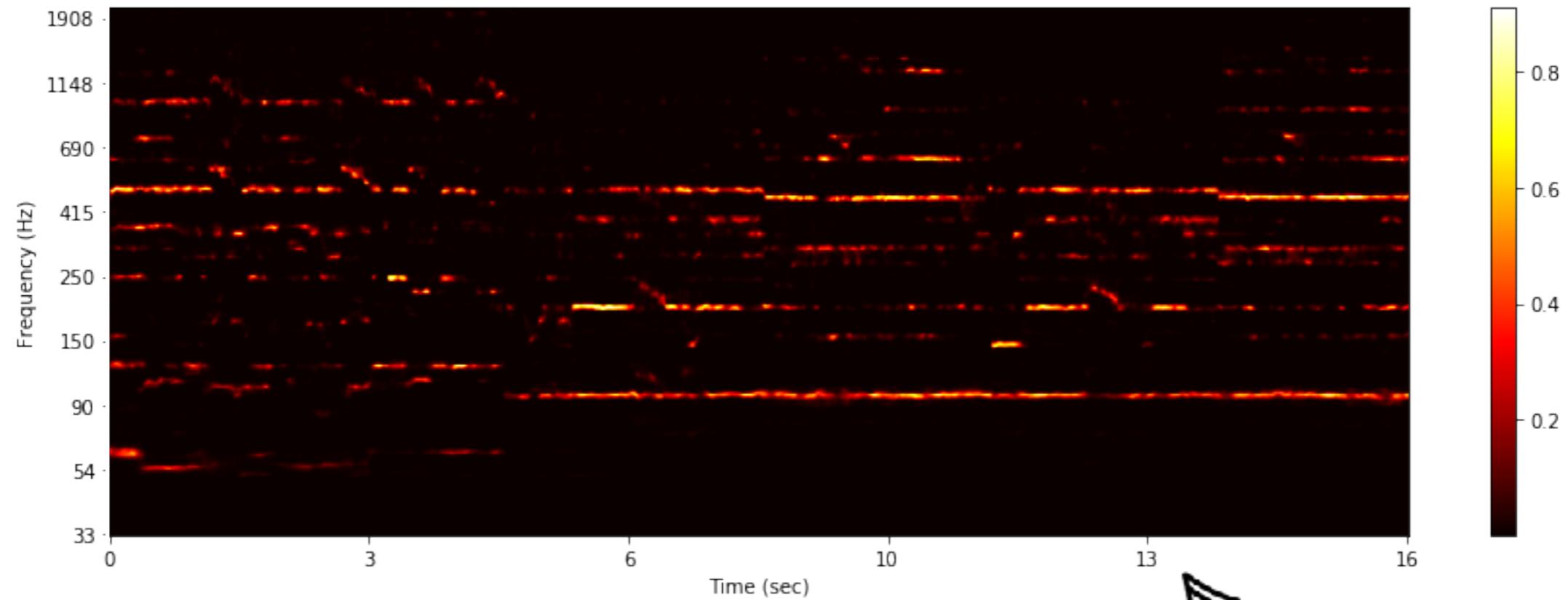
```
('Precision', 0.881),  
('Recall', 0.886),  
('Accuracy', 0.791),  
('Substitution Error', 0.047),  
('Miss Error', 0.067),  
('False AlarmaError', 0.072),  
('Total Error', 0.187),  
('Chroma Precision', 0.900),  
('Chroma Recall', 0.905),  
('Chroma Accuracy', 0.823),  
('Chroma Substitution Error', 0.028),  
('Chroma Miss Error', 0.067),  
('Chroma False Alarm Error', 0.072),  
('Chroma Total Error', 0.167)
```



Frequency (Hz)

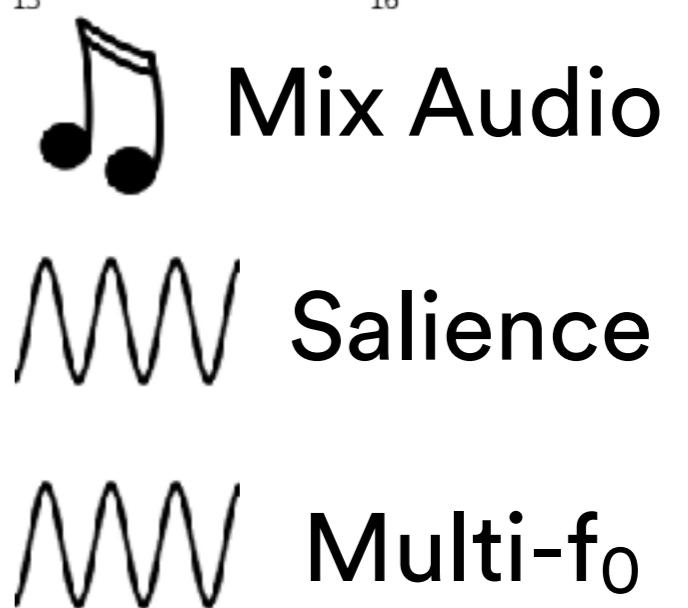
[github.com/craffel/mir\\_eval](https://github.com/craffel/mir_eval)  
mir\_eval.multipitch.evaluate

# Saliency Qualitative Evaluation



$$y(t) = a(t) \times \sin(2\pi \int^t f(x) dx)$$

- >  $a(t)$  - “amplitude” over time  
voicing → binary
- “confidence” → continuous
- >  $f(t)$  - frequency estimate



# Datasets

# f<sub>0</sub> Datasets



	<b>Year</b>	<b>Dur.</b>	<b>Multi-f0</b>	<b>Melody</b>	<b>Vocal</b>	<b>Bass</b>	<b>Piano</b>	<b>Guitar</b>
<b>RWC</b>	2002	7 h		notes				
<b>MAPS</b>	2008	18 h	notes				notes	
<b>MIR-1K</b>	2009	3 h		notes				
<b>Bach10</b>	2012	5 m	notes					
<b>MedleyDB</b>	2014	7 h		f0	f0	f0 (1h)		
<b>iKala</b>	2015	2 h		f0	f0			
<b>Su-AMT</b>	2015/6	12 m	notes					
<b>Weimar Jazz</b>	2016	13 h		notes		notes (1h)		
<b>OrchSet</b>	2016	23 m		f0				
<b>Lakh MIDI</b>	2016	15 d	notes					
<b>MusicNet</b>	2018	34 h	notes					
<b>GuitarSet</b>	2018	3 h	notes					notes

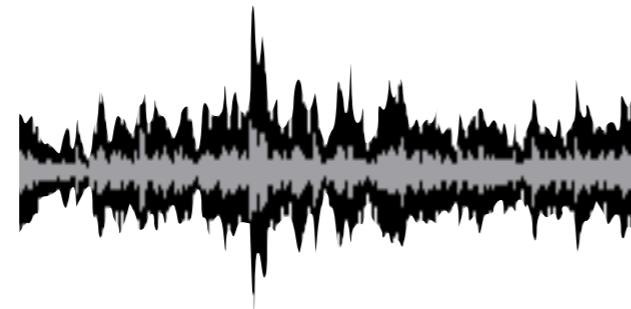
# f<sub>0</sub> Datasets

	Link
RWC	<a href="https://staff.aist.go.jp/m.goto/RWC-MDB/rwc-mdb-p.html">https://staff.aist.go.jp/m.goto/RWC-MDB/rwc-mdb-p.html</a>
MAPS	<a href="http://www.tsi.telecom-paristech.fr/ao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/">http://www.tsi.telecom-paristech.fr/ao/en/2010/07/08/maps-database-a-piano-database-for-multipitch-estimation-and-automatic-transcription-of-music/</a>
MIR-1K	<a href="https://sites.google.com/site/unvoicedsoundseparation/mir-1k">https://sites.google.com/site/unvoicedsoundseparation/mir-1k</a>
Bach10	<a href="http://music.cs.northwestern.edu/data/Bach10.html">http://music.cs.northwestern.edu/data/Bach10.html</a>
MedleyDB	<a href="http://medleydb.weebly.com/">http://medleydb.weebly.com/</a>
iKala	<a href="http://mac.citi.sinica.edu.tw/ikala/">http://mac.citi.sinica.edu.tw/ikala/</a>
Su-AMT	<a href="https://sites.google.com/site/lisupage/research/new-methodology-of-building-polyphonic-datasets-for-amt">https://sites.google.com/site/lisupage/research/new-methodology-of-building-polyphonic-datasets-for-amt</a>
Weimar Jazz	<a href="https://jazzomat.hfm-weimar.de/dbformat/dboverview.html">https://jazzomat.hfm-weimar.de/dbformat/dboverview.html</a>
OrchSet	<a href="https://www.upf.edu/web/mtg/orchset">https://www.upf.edu/web/mtg/orchset</a>
Lakh MIDI	<a href="http://colinraffel.com/projects/lmd/">http://colinraffel.com/projects/lmd/</a>
MusicNet	<a href="https://homes.cs.washington.edu/~thickstn/musicnet.html">https://homes.cs.washington.edu/~thickstn/musicnet.html</a>
GuitarSet	<a href="https://guitarset.weebly.com/">https://guitarset.weebly.com/</a>

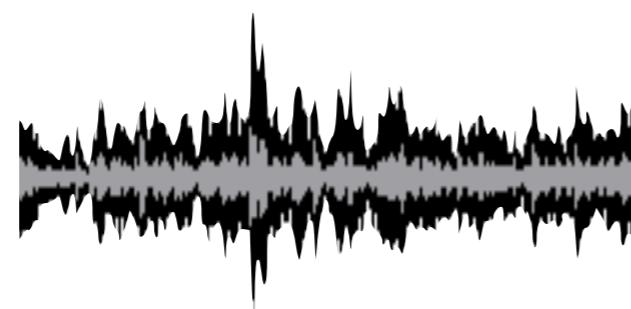
# Data Generation

# How do we obtain $f_0$ labels?

Polyphonic Audio

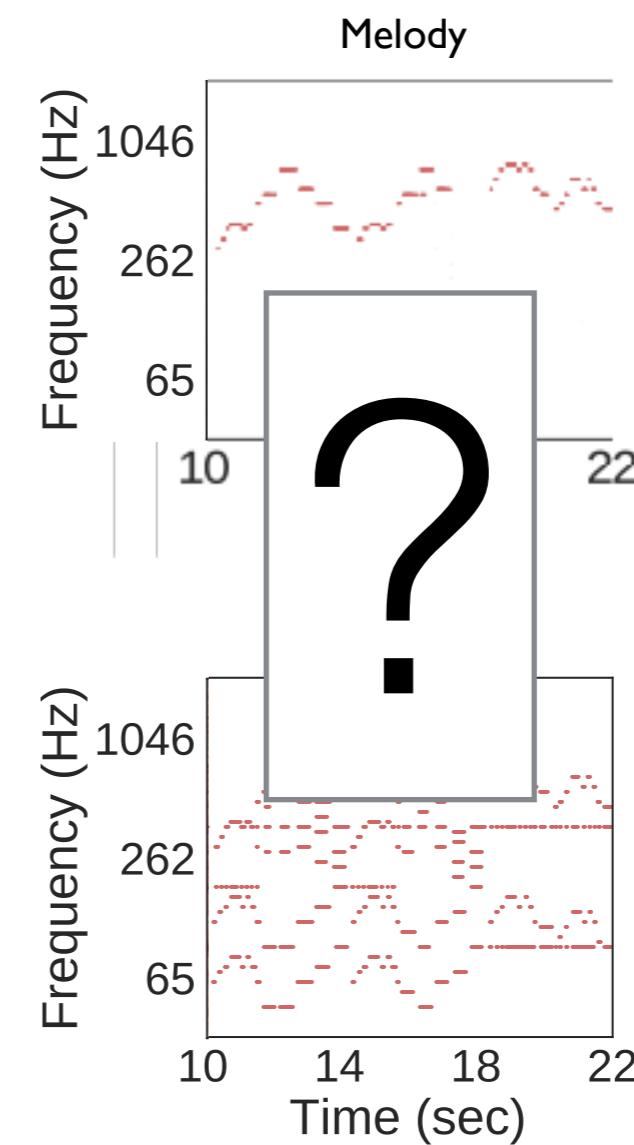


+

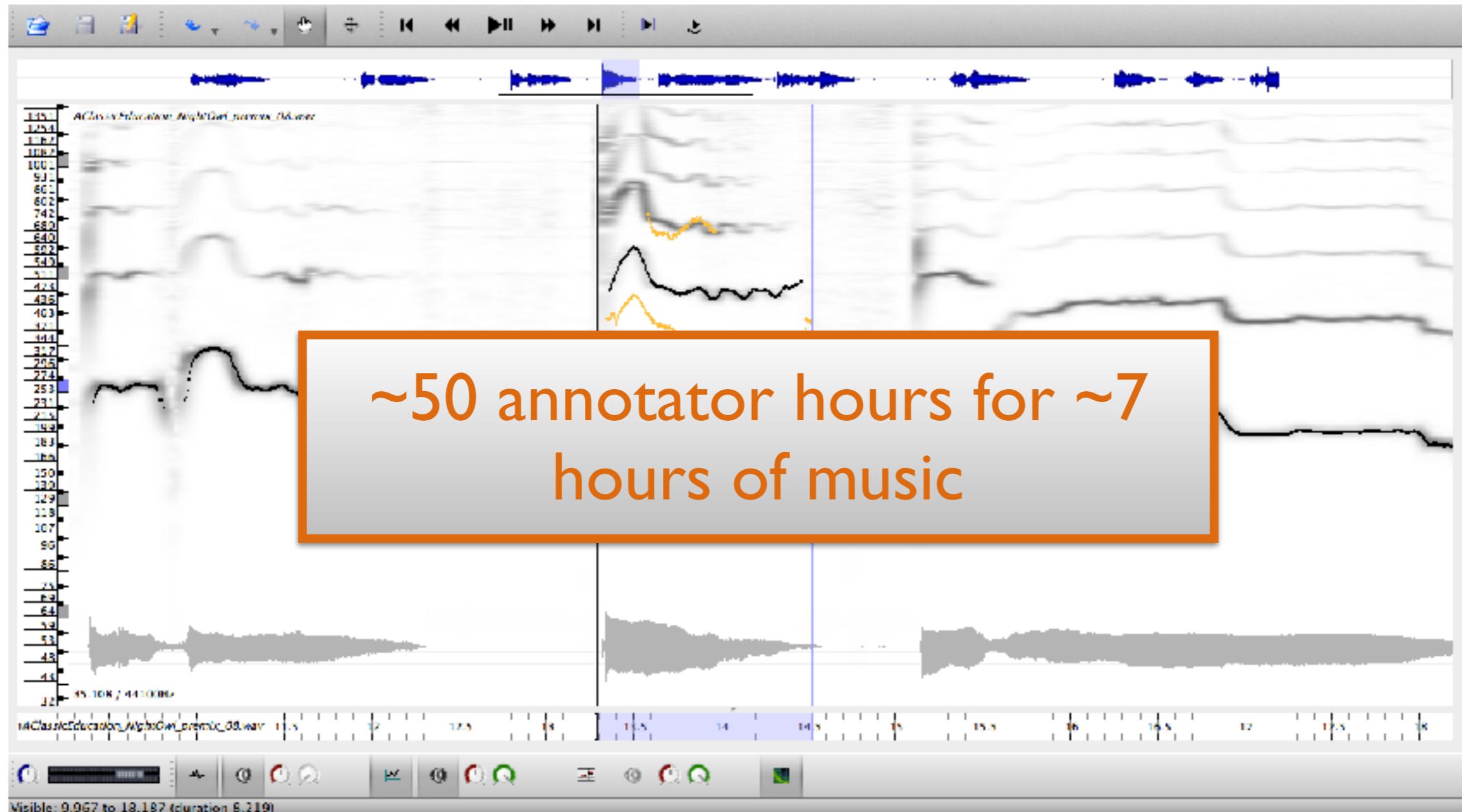


+

$f_0$  Annotations

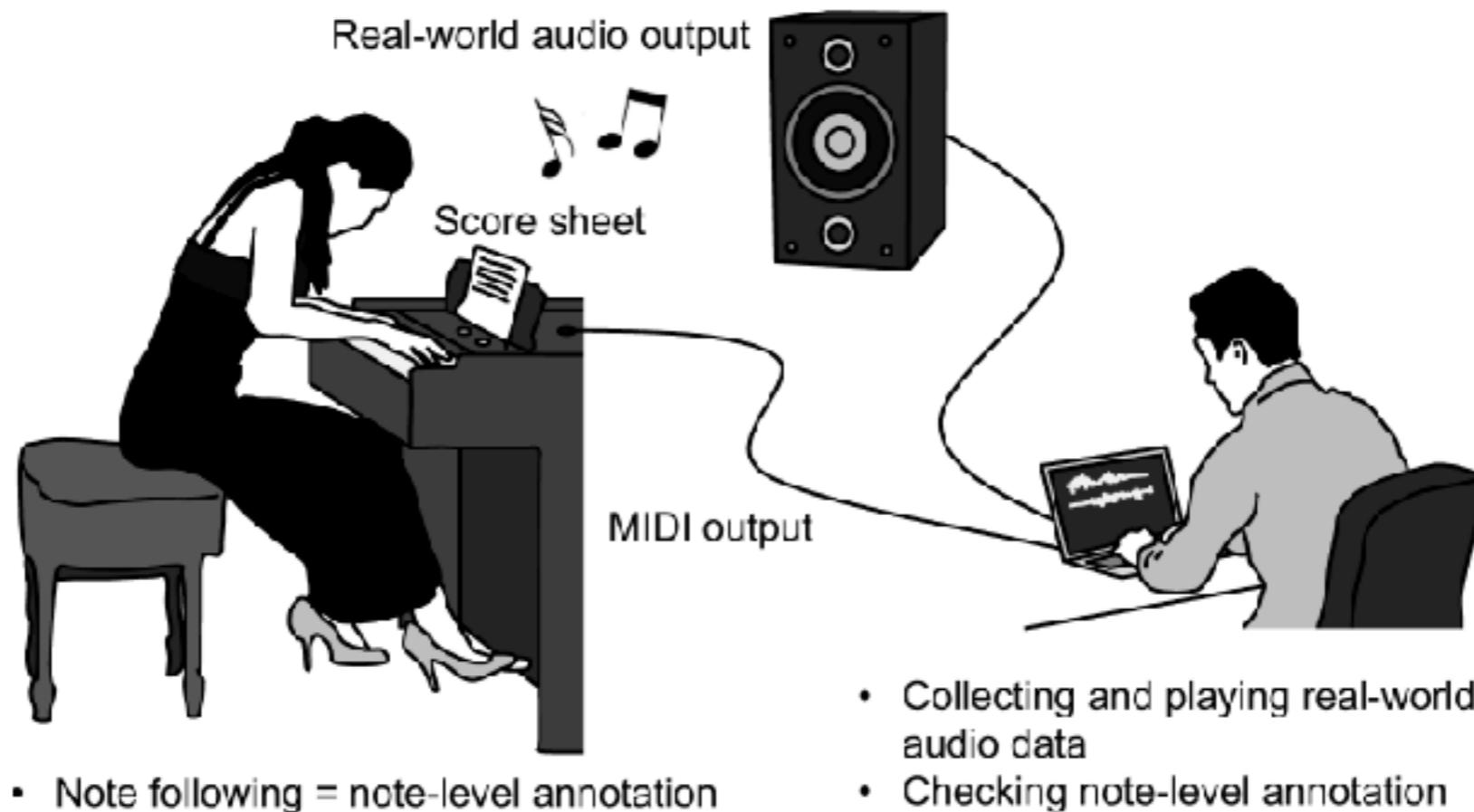


# Manual f<sub>0</sub> correction



<https://code.soundsoftware.ac.uk/projects/tony>

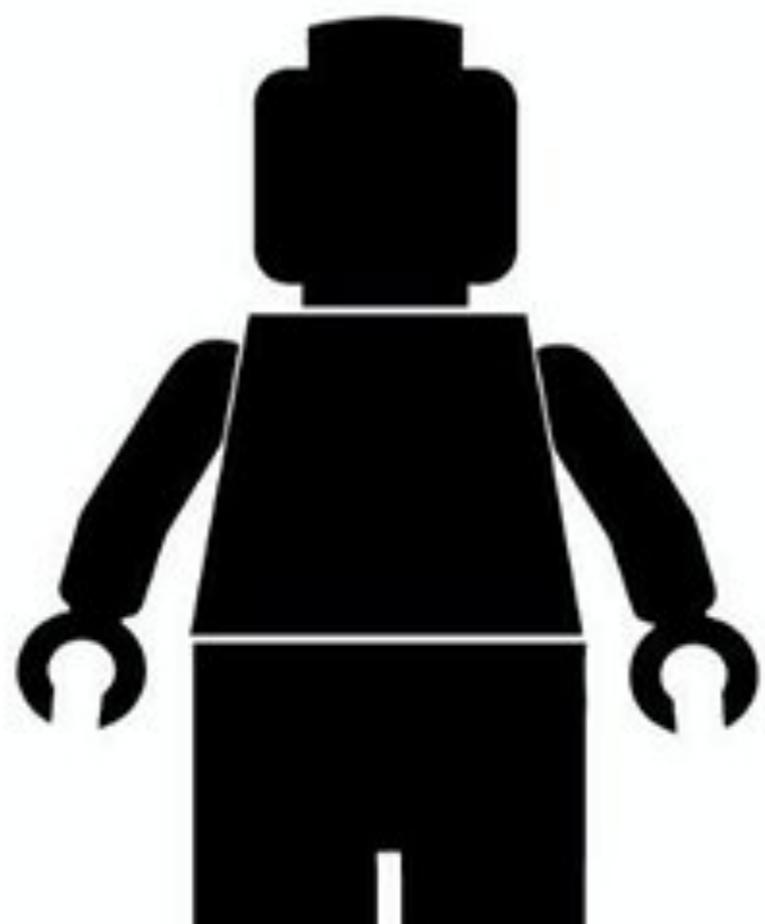
# Assisted Note Annotations



Li Su, Yi-Hsuan Yang

"Escaping from the Abyss of Manual Annotation: New Methodology of Building Polyphonic Datasets for Automatic Music Transcription"  
ICMC (2015)

# Synthesized MIDI



# Annotation by Resynthesis



Played



Estimated f0

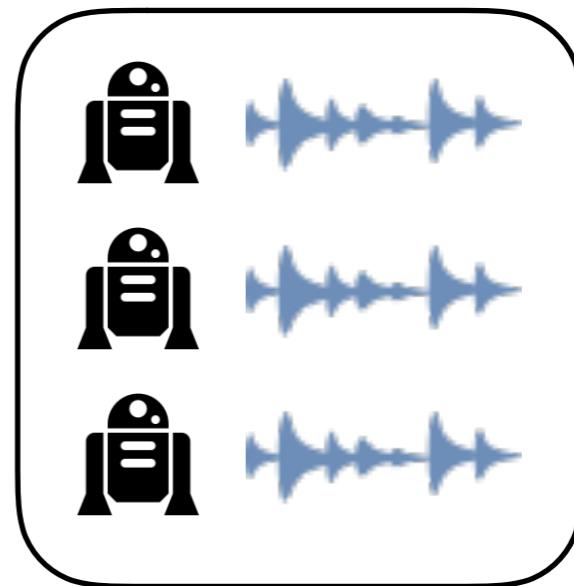
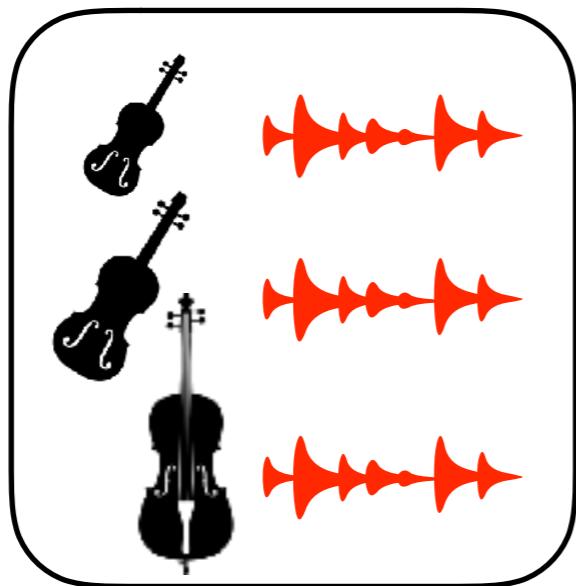


Synthesized f0



Justin Salamon, Rachel Bittner, Jordi Bonada, Juan José Bosch,  
Emilia Gómez, Juan Pablo Bello  
“An Analysis/Synthesis Framework for Automatic f0 Annotation of  
Multitrack Datasets”  
ISMIR (2016)

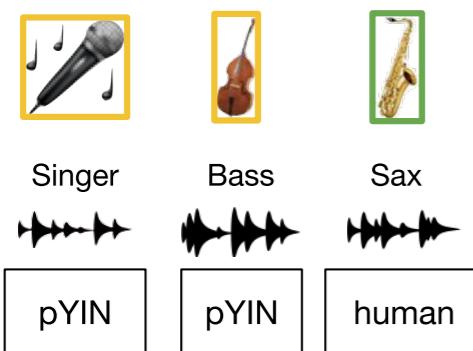
# Annotation by Resynthesis



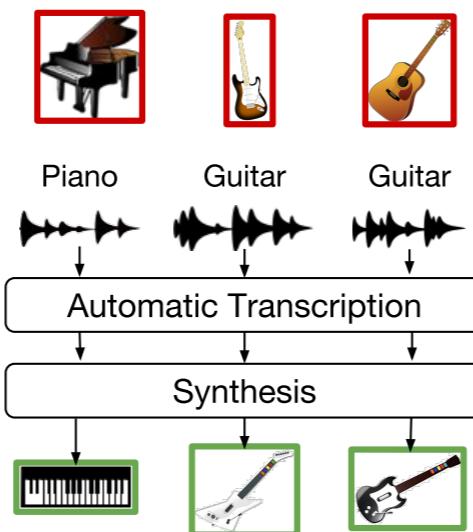
Justin Salamon, Rachel Bittner, Jordi Bonada, Juan José Bosch,  
Emilia Gómez, Juan Pablo Bello  
“An Analysis/Synthesis Framework for Automatic f0 Annotation of  
Multitrack Datasets”  
ISMIR (2016)

# Automatic Multitrack Annotations

## Available Instrument Annotations



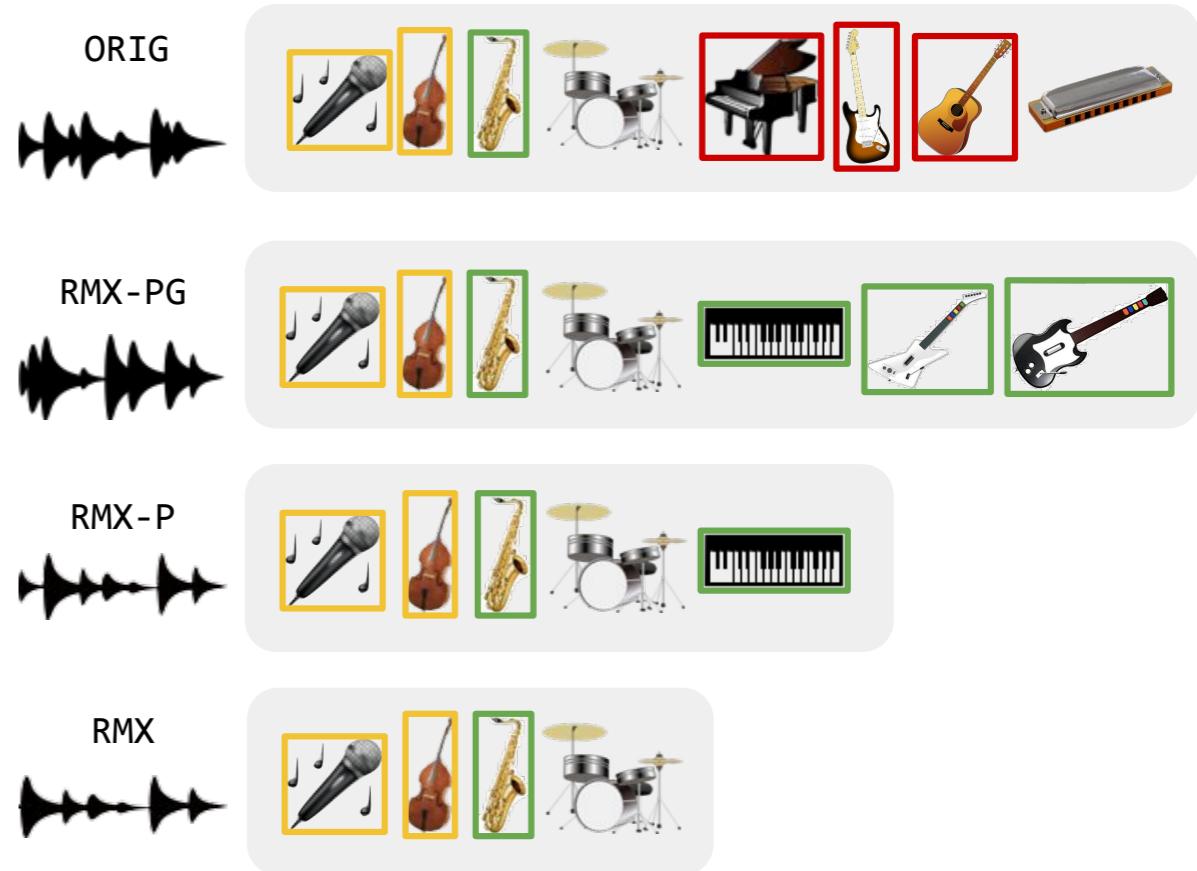
## Instrument Annotations via Synthesis



## Individual Instrument Annotation Quality



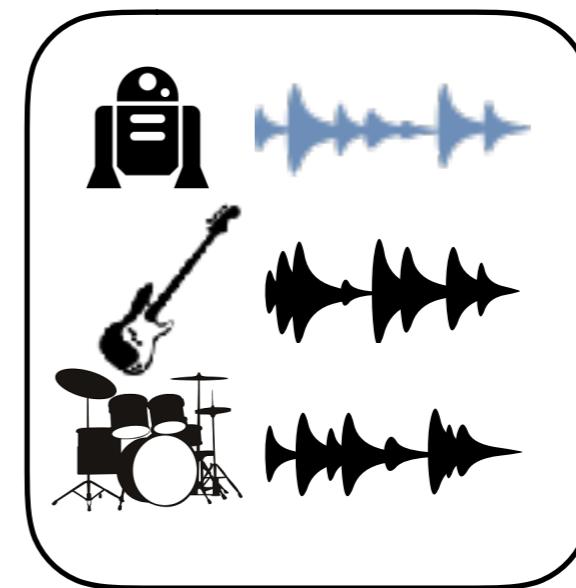
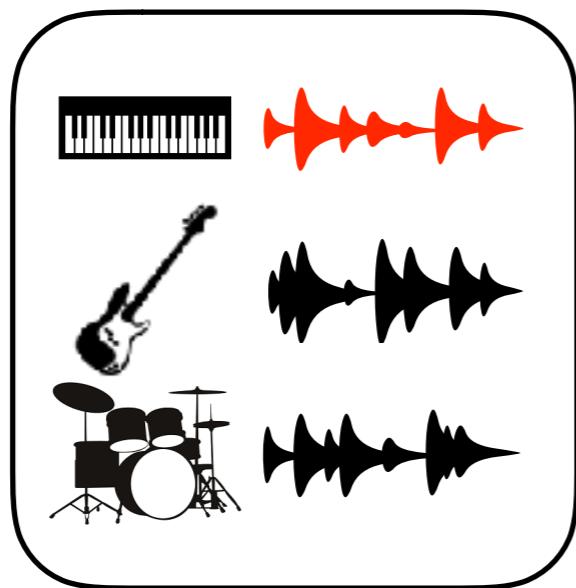
## Mixture Types and Available $f_0$ Annotations



## Mixture Annotation Types (this example)

	ORIG	RMX-PG	RMX-P	RMX
Multiple- $f_0$	x	✓	✓	✓
Melody	✓	✓	✓	✓
Bass	✓	✓	✓	✓
Vocal	✓	✓	✓	✓
Piano	x	✓	✓	✓=0
Guitar	x	✓	✓=0	✓=0

# Piano/Guitar Replacement Examples



# Summary

- Most polyphonic f0 estimation models produce a pitch salience representation and decode to f0 or notes
  - Tasks are approached separately most of the time
  - Transcription and f0 estimation are almost the same
- Evaluation metrics are primarily frame based
- Datasets for each individual task are limited
- Multitracks can be used to generate data more easily

# Open Problems

- Perceptually motivated evaluation metrics
  - go beyond frame level
  - incorporate salience
- Optimization function that is an upper bound for f0 evaluation metrics
- Directly predict outputs: f0 time series or notes
- Joint source separation + f0 estimation
- Include onset information for f0 trajectories
- Train a simple system on a ton of data and see what happens