

The image displays a musical score for piano, featuring three systems of staves. The first system includes dynamic markings such as *cresc.*, *f*, and *p*, along with chord annotations like V , $vii^{\circ}4$, I^6 , iv^6 , $V^{(4)}/V$, V/V , iv^6_4 , V^2 , I^6 , $V(\sharp)$, and V^{\sharp}_5 . The second system shows *fp*, *p*, and *f* dynamics, with chords $I\{$, $vii^{\circ}6$, I , IV , V , and I^6 . The third system includes *fp*, *p*, and *f* dynamics, with chords V/V , V^{\sharp}_5 , I , $vii^{\circ}6$, I^6 , IV , $V(\sharp)$, $I(\sharp)$, and I . A red box on the right contains the text "Aligning audio to annotated score labels". A dark grey box in the lower center contains the text "Clémentine Lévy-Fidel", "Life Sciences Engineering", and "Semester project".

Aligning audio to annotated score labels

Clémentine Lévy-Fidel
Life Sciences Engineering
Semester project

- Motivation
- Applications
- Related work
- Approaches
 - Beat tracking and meter reconstruction
 - Chord detection
 - Score following
 - Dynamic Time Warping
- The Aligner tool
- Next steps
- Project highlights
- Acknowledgments

Motivation – score labels to audio alignment

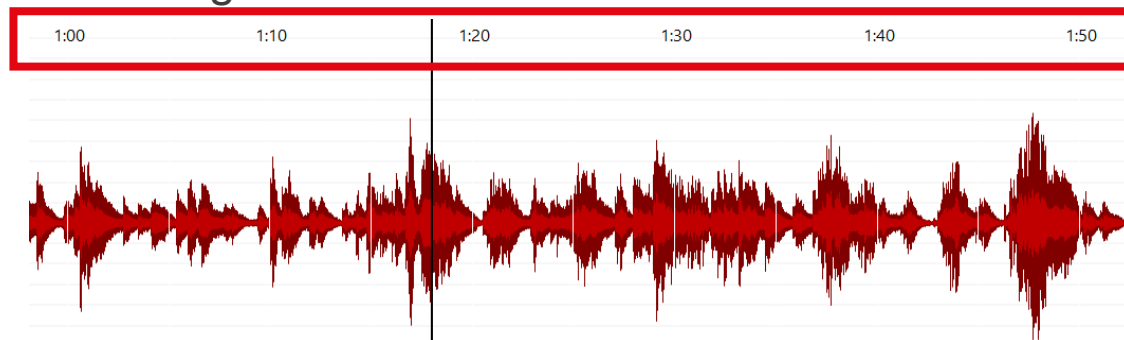
■ Score labels

A musical score snippet in G major, measures 9-12. The score includes dynamic markings (*fp*, *p*, *f*) and articulation markings (*tr*). A red box highlights the harmonic analysis labels below the staff:

p *fp* *p* *fp* *f* *p*

I⁶ vii^{o6} I⁶ IV⁶ V⁶ vii^{o6} I⁶ V⁶/ii ii⁽⁴⁾ vii^o I⁶ {

■ Audio recording



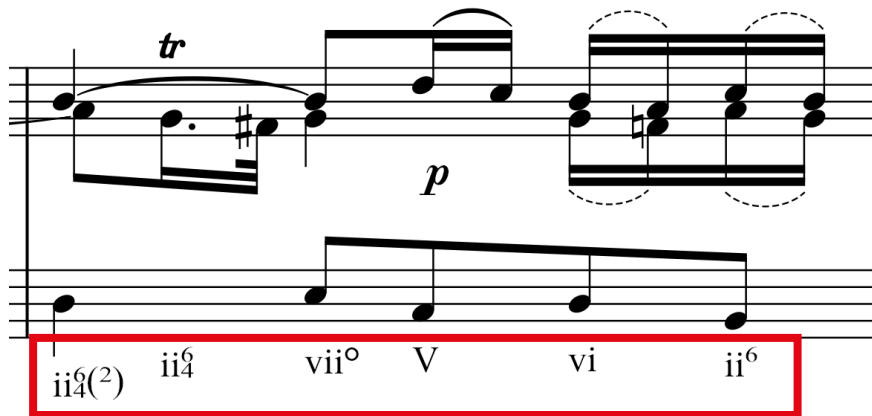
Motivation – score labels to audio alignment

- **Goal:** automate alignment
- **Desired outcome:** temporal positions of labels

<i>Timestamp</i>	<i>Label</i>
0.90	F.I{
1.48	viiio6
2.22	I
4.96	IV
5.66	V
...	...

Challenges:

- Extensive count of labels per piece:
 - Annotated Mozart Sonatas corpus: 104 to 756 labels
 - Time-consuming manual alignment



A musical score snippet showing two staves. The top staff contains a trill (tr) and a piano (p) dynamic marking. The bottom staff contains a series of notes. Below the bottom staff, a red box highlights a sequence of harmonic labels: $ii_4^6(2)$, ii_4^6 , vii^o , V , vi , and ii^6 .

Motivation – score labels to audio alignment

Challenges:

- Varying annotation or audio profiles

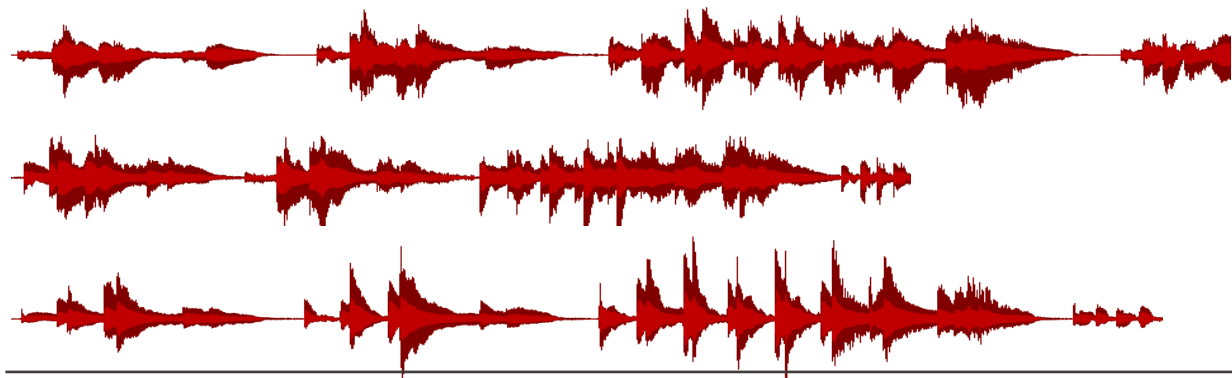
A musical score for piano in 3/4 time. The score consists of four measures. The first measure has a piano (*p*) dynamic and a fortissimo (*fp*) dynamic. The second measure has a piano (*p*) dynamic and a fortissimo (*fp*) dynamic. The third measure has a fortissimo (*f*) dynamic. The fourth measure has a fortissimo (*f*) dynamic and a piano (*p*) dynamic. The score includes various articulations such as slurs, trills (*tr*), and accents. Below the notes, there are Roman numerals and chord symbols: F.II, vii⁰⁶, IV, V, I, V⁶, I, f, vi⁶, ii⁽⁴⁾, V⁷, I, and {.

Performer

Fazil Say

Marta Deyanova

Roberte Mamou



- Practical usability:
 - Transfer annotations from one recording to another
 - Switch between music formats: audio, score, MIDI...
- Behavioural experiments:
 - Align with EEG measurements
 - Align with assessments of perceived musical tension
- Support music theory:
 - Understand experts' annotations live
 - Live visualisation in context of presentations, teaching...
- Extend the use of the *Annotated Mozart Sonatas* corpus
 - 18 sonatas, 56 scores
- Apply to other databases and research, e.g. *Montreux Jazz Festival*

- **Meter inference** – via stochastic models

Beat and downbeat tracking systems:

- Bar pointer model (Whiteley et. al, 2006)

- **Music tracking** – via stochastic models

Score following for live performance accompaniment:

- *AnteScofo* (IRCAM, 2009)
- *MusicPlusOne* (Christopher Raphael, 2001)
- Deep-learning based approaches on music sheet images

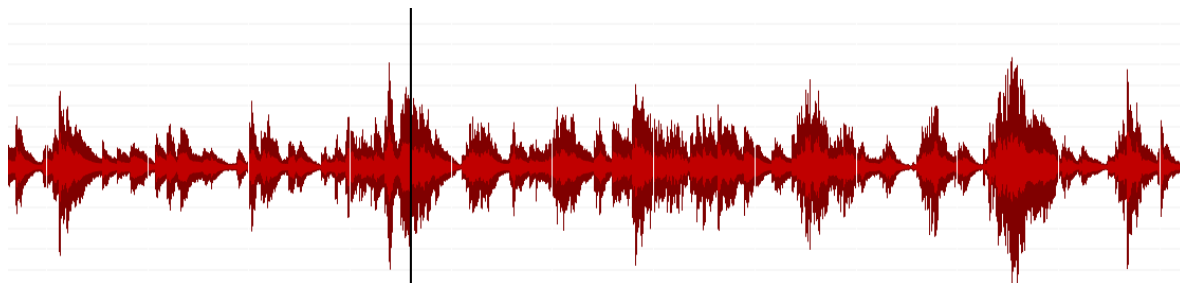
- **Music synchronization** – via Dynamic Time Warping

International Audio Laboratories Erlangen:

- Fundamentals of Music Processing (FMP) (Müller et. al, 2015)
- *LibFMP* (Müller et. al, 2020)
- *SyncToolBox* (Müller et. al, 2021)

Approaches – 1. Beat tracking

Detection by beat [1] and downbeat [2] tracking algorithms from madmom's models:



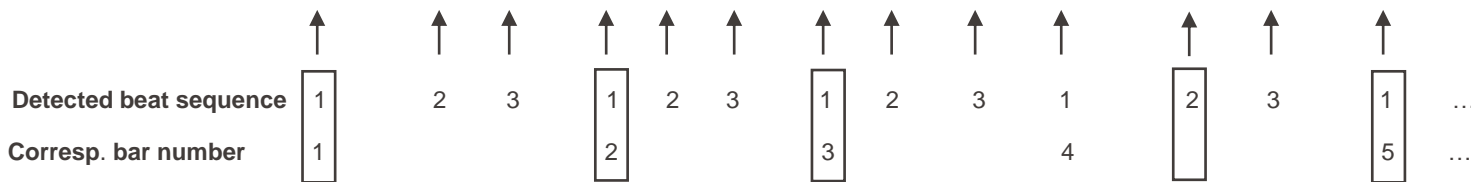
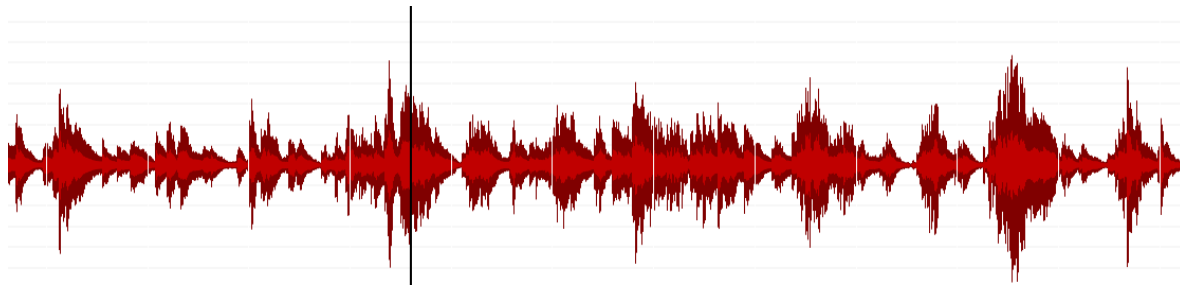
Detected beat sequence	1	2	3	1	2	3	1	2	3	1	2	3	1	...
Corresp. bar number	1			2			3			4			5	...

[1] Krebs et. al, 2015

[2] Böck et. al, 2016

Approaches – 1. Beat tracking

Alignment:

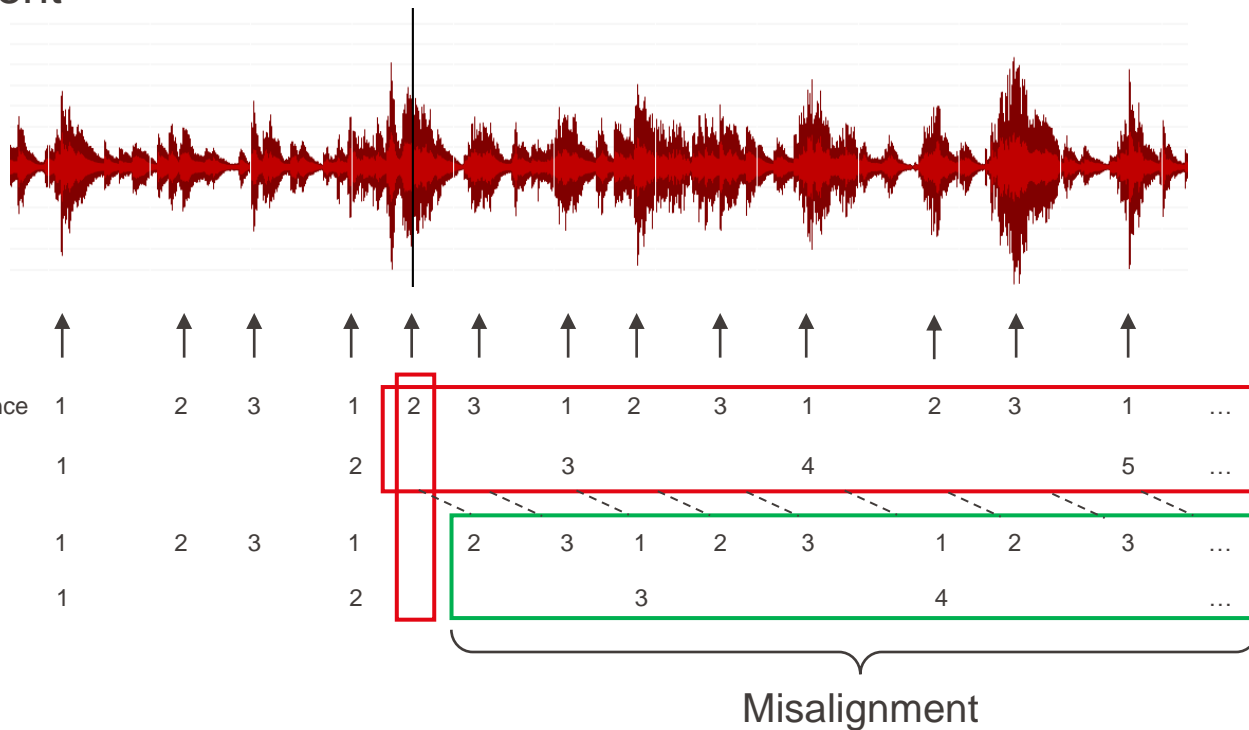


+

Beat	1	1	2	1	2	...
Bar	1	2	2	3	4	...
Label	C	Dm7	G7	C	Am	...

Pitfalls:

■ Misalignment

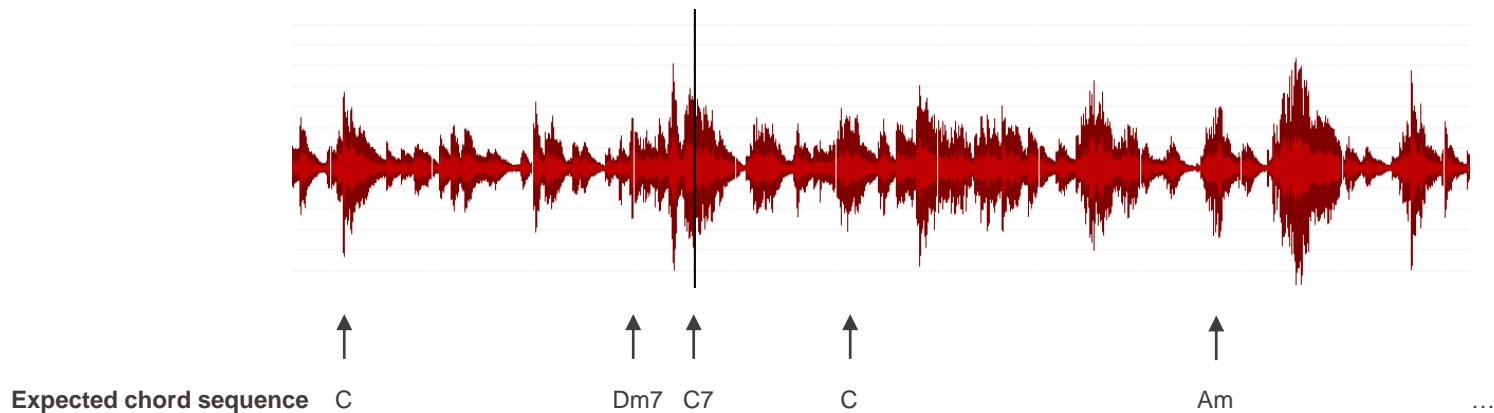


Pitfalls:

- Misalignment
- Phase locks with a correct metrical grid but not necessarily at the intended beat level:
 - Interprets half beats as beats
 - In general, happens to detect a multiple of the number of beats (“octave error”)

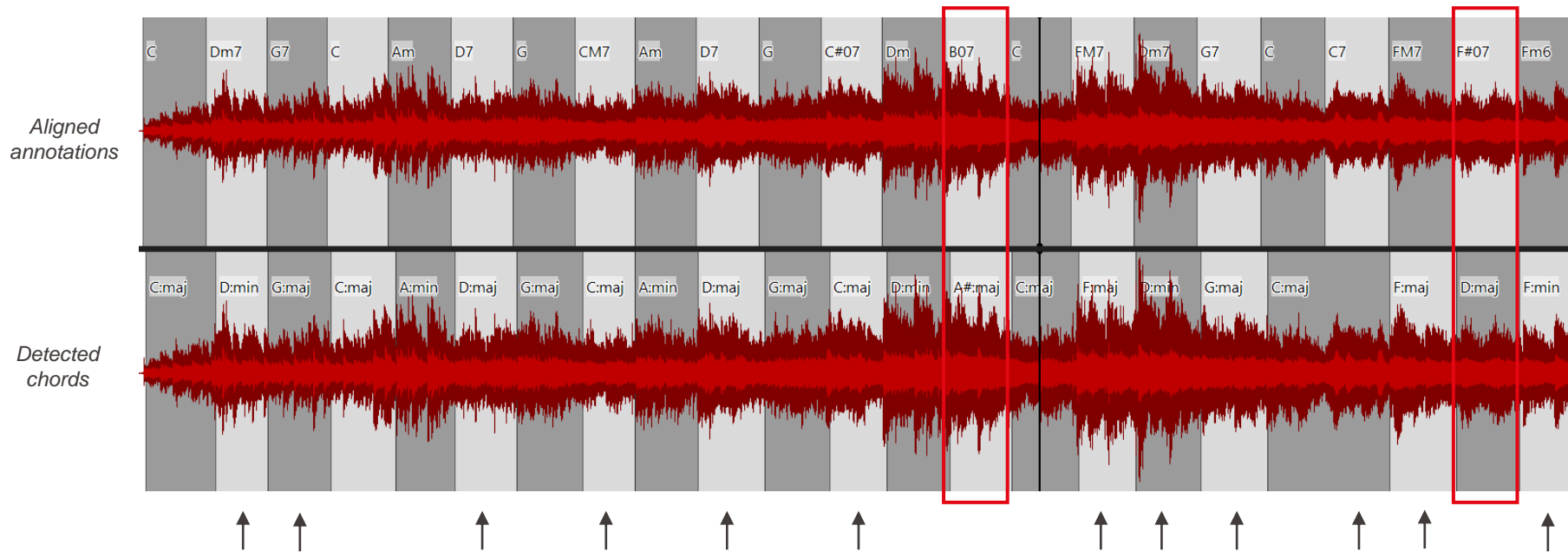
Approaches – 2. Chord detection

Chord detection by chroma detection algorithms from madmom's models [3]:



[3] Korzeniowski et. al, 2016

Bach, Prelude in C



Approaches – 2. Chord detection

Pitfalls:

- Not robust enough:
 - Simplified chords (only major or minor)
 - Frequent false alarms / misses
- Correctness can vary upon annotator

The image shows a musical score snippet. The top staff is in treble clef, marked 'pp' (pianissimo) and 'calando' (rushing). The bottom staff is in bass clef. A bracket under the bass line indicates a detected chord.

Annotator: C64
Detected chord: F

Approaches – 3. Score following

Deep learning approach reading score images [4]

Pitfalls:

- More demanding to use score images than XML or MIDI files if they are provided
- Method difficult to adapt
- Accuracy not entirely guaranteed

[4] Henkel et. al, 2021

Possible workarounds for misalignment pitfalls:

- Use state-of-the-art sequence alignment algorithms
- Train *madmom*'s models or score following models for our type of music
 - higher correct detection rate
 - more robust alignment

Remaining problems:

- Need of manual verification
- Lack of aligned data for comparison with ground truth

Synchronization approach:

For every event in the score sequence (i.e., note), match with a corresponding event in the audio sequence

→ Dynamic Time Warping:

Find alignment path between two time series by minimizing distance

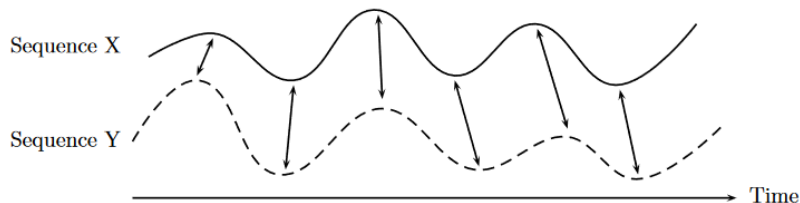


Fig. 4.1. Time alignment of two time-dependent sequences. Aligned points are indicated by the *arrows*

Image source: Müller et al., 2007.

Following Ewert et al. (2009):

- Audio feature extraction

- Chroma onset (CO) features combining:

- Chroma filtering
- Onset energy peak picking

- Locally adaptive normalization

→ normalize with local maximum of sequence to make CO features invariant to dynamic variations

- Add temporal decay

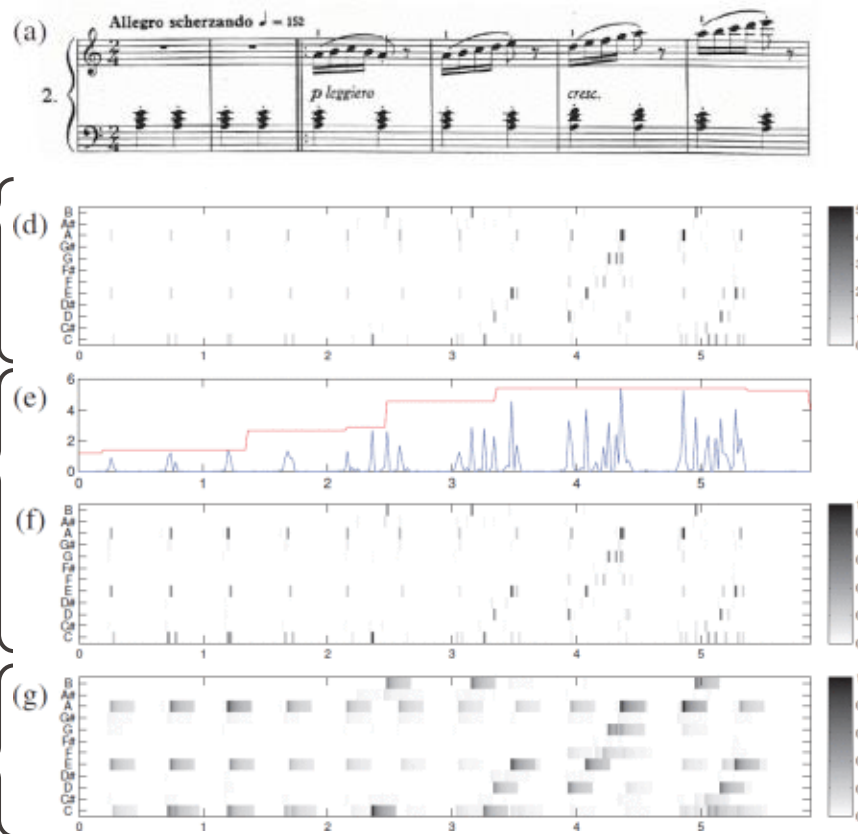


Fig. 1. (a) First six measures of Burgmüller, Op. 100, Etude No. 2. (b) - (g) feature representations of a corresponding audio recording (see Sect. 2 for a description).

Following Ewert et al. (2009)

- Synchronization algorithm
 - Similarity (or cost) matrix: evaluates local similarity cost measure for each pair of features between both sequences
 - Determine optimum-cost alignment path

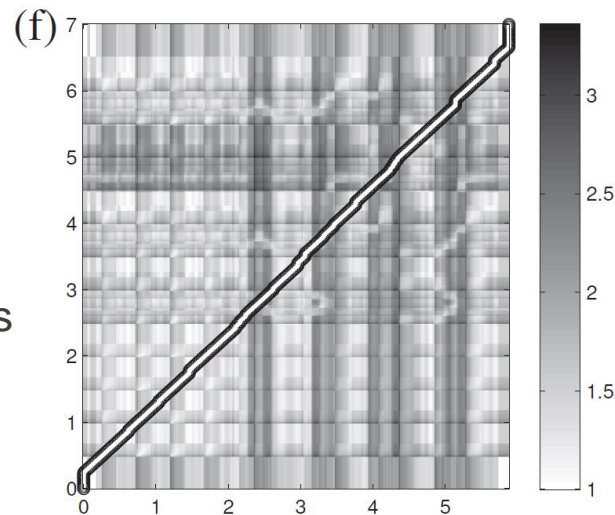


Fig. 2. (a)-(c) Illustration of the effect of the decay operation on the cost matrix level. (d) C_{chroma} , (e) C_{DLNCO} (f) $C_{\text{chroma}} + C_{\text{DLNCO}}$ for Burg2.

Applying DTW to the Annotated Mozart Sonatas corpus

Dataset content:

■ Notes

quarterbeats	duration_qb	mc	mn	mc_onset	mn_onset	timesig	staff	voice	duration	gracenote	nominal_duration	scalar	tied	tpc	midi	chord_id
0	0.375	1	1	0	0	3/4	1	1	3/32		1/16	3/2		3	69	0
3/8	0.125	1	1	3/32	3/32	3/4	1	1	1/32		1/32	1		-1	65	1
1/2	0.375	1	1	1/8	1/8	3/4	1	1	3/32		1/16	3/2		-2	70	2

■ Harmonies

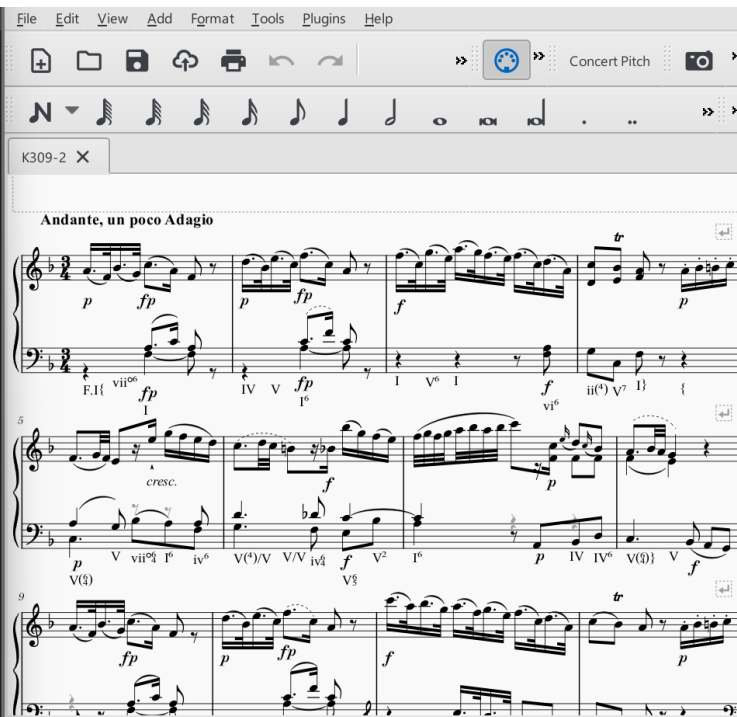
quarterbeats	duration_qb	mc	mn	mc_onset	mn_onset	timesig	staff	voice	label	globalkey	localkey	pedal	chord	numeral	form	figbass
0	0.5	1	1	0	0	3/4	2	1	F.I{	F	I		I	I		
1/2	0.5	1	1	1/8	1/8	3/4	2	1	viiio6	F	I		viiio6	vii	o	6
1	2.0	1	1	1/4	1/4	3/4	2	1	I	F	I		I	I		

Input: notes, labels, audio

Steps:

- Work on **notes** events: convert **symbolic annotations** to **synthesized audio** representation
 - Temporal occurrence, duration, pitch (midi), velocity, instrument
 - Extract similar features as out of audio
- Perform **synchronization to audio** by DTW
- Join **notes** and **labels** datasets on *quarterbeats*

Output



Github repository (<https://github.com/clelf/Aligning-audio-to-annotated-score-labels>) containing:

- Dataset preparation tutorial (using `ms3` parser)
- Code to adapt *SyncToolBox* to the use of the *Mozart Sonatas Annotated* dataset
- Command line interface:

```
python aligner.py -a [audio_WAV_file] -n [notes_TSV_file]  
-l [labels_TSV_file] -o [CSV_file_to_write_results_to]
```


Additional features

- Output aligned notes for score following purpose
- Evaluation module to be used if ground truth data is labelled

- Generalize pipeline to data outside the *Mozart Sonatas* corpus
- Generalize algorithm to scores and/or audio recordings that contain repetitions, following Grachten et. al (2013)
- Label data manually and use ML-based approaches
- Investigate JKU's *matchmaker* library when published

- Invaluable help of Steffen, Gabriele and Johannes 😊
- Great introduction to the field of audio processing, music information retrieval, digital musicology ...
- Contact with Gerhard Widmer's lab at Johannes Kepler University (Linz)
- Creating a useful tool!

Special thanks to Steffen, Gabriele and Johannes for their support

Thanks to the whole of DCML team for the warm welcome and insightful conversations



Thank you!

Questions?