# Automated Music Transcription based on Formal Language Models

Florent Jacquemard, Inria Paris

Philippe Rigaux
le cnam Paris

Florent Jacquemard
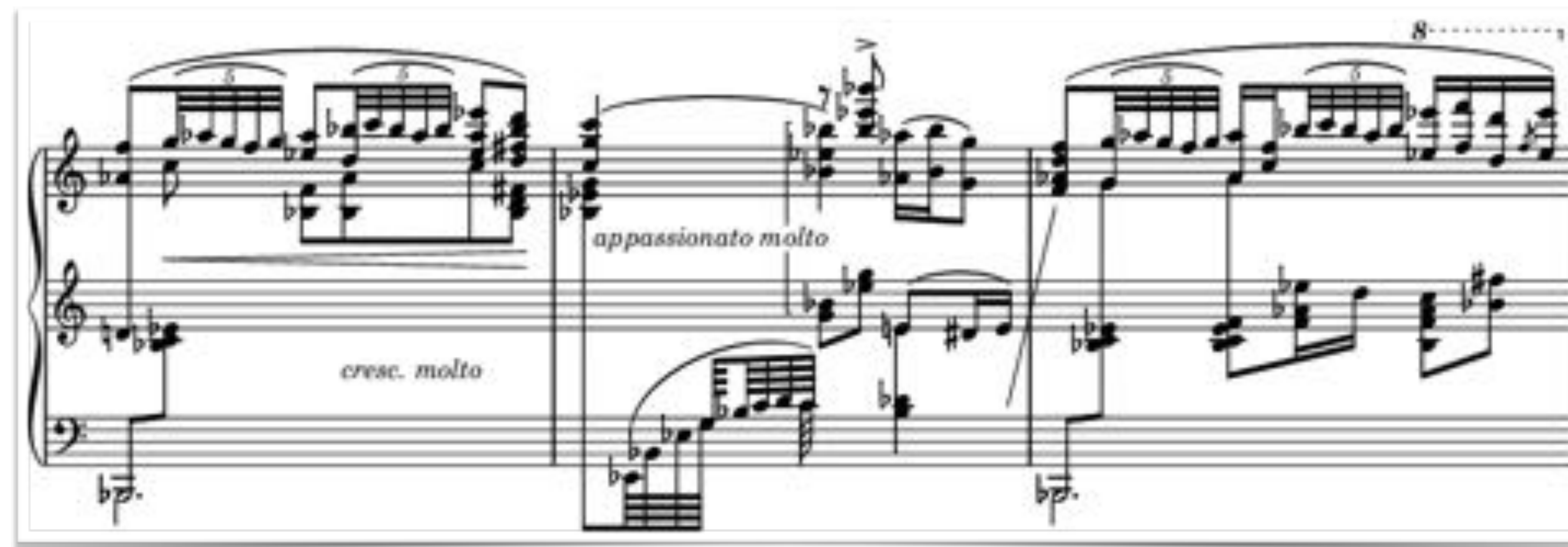Inria informatiques mathématiques

Raphaël Fournier-S'niehotta
le cnam

Lydia Rodriguez-de la Nava
PhD (Codex, Inria)

Tiange Zhu
PhD (Polifonia, H2020)

post-doc (Collabscore, ANR)

## Music Notation Processing



E. Granados, Goyescas
typesetted with Lilypond

Western Music Notation = graphical format for music practice,
in use since ~1000 years (Guido d'Arezzo)



vs



Philippe Manoury
Tensio for string quartet and electronics

(digital) music scores, a natural language  for

• performers
  performance : real-time reading or memoization

• composers
  authoring, exchange

• teachers & students
  transmission

• editors
  access digital score libraries e.g. nkoda.com

• librarians
  cultural heritage preservation: e.g. Gallica

• scholars (historians, musicologists…)
  research, analysis

Philippe Rigaux
le cnam Paris

Florent Jacquemard
*Inria* informatiques mathématiques

Raphaël Fournier-S'niehotta
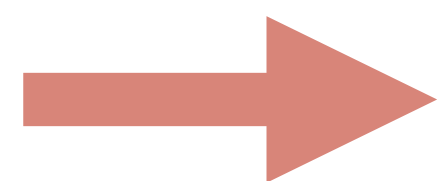le cnam

Lydia Rodriguez-de la Nava
PhD (Codex, Inria)

Tiange Zhu
PhD (Polifonia, H2020)

post-doc (Collabscore, ANR)

## Music Notation Processing

- Structured music score models
  hierarchical representation of music scores
- Music scores languages
  finite representations of ∞ sets of scores (*style*)
- Search and retrieval
  indexing, exact or approximate search, facetted
- Similarity metrics
  string and tree edit-distances

## Applications

- Databases of digital music scores
  Cultural heritage preservation  H2020 Polifonia - U. Bologna,
  Open University, King's College, Vrije U. Amsterdam
- Computational Musicology
  neuma.huma-num.fr - IReMus (Paris), AlgoMus (Lille)
- Optical Music Recognition, Crowdsourced correction
  ANR Collabscore - IRISA (Renne), French National Library, Royaumont
- Automated Music Transcription
  JSPS 採譜, grant Yamaha Music Foundation - JAIST, Nagoya U.

Conversion of a recorded music performance into a music score ~ *speech-to-text* in NLP
a holy graal in Computer Music since 1970's



©1976 **Nature Publishing Group**

Conversion of a recorded music performance into a music score

**source(s)**

Audio recording

MIDI device
(score edition)

Algorithmic composition
DAW

audio Music Information Retrieval
- fundamental freq. estimation
- onset detection
- beat tracking …

**intermediate representation**
**piano roll** (MIDI file)
- unquantized onsets, durations
- quantized pitches

symbolic Music Information Retrieval
- rhythm quantization
- tempo tracking
- score engraving…

**target**
**music score**
(XML file)

Rhythm quantization with grids, *e.g.* MIDI files import
• in score editors (Finale, Sibelius, Dorico, Musescore…),
• or in DAWs (Ableton Live, Logic…)

Alignment of every input time point (onset) to the closest position
in a *grid* = sequence of equidistant time position.



input

grid 16th note

grid 32th note

hierarchical grid

alignment

poor fit, good readability

good fit, bad readability

closer to intuition

## regular grids

- search of a best quantization is possible by a brute-force enumeration: 8th note grid, 16th, 32th, 64th…
- result not always optimal
- problems with tuplets (so called *"irrationals"* 3, 5, 7…*)*

## hierarchical grids

- more "natural" results
- brute force enumeration impossible
- how to specify the grids to try ?

beamed

unbeamed

hierarchical
note
durations

Polonaise in D minor from Notebook for Anna Magdalena Bach    BWV Anh II 128

**metric structure**

bar    1                              2            3                    4                        5
beat   1.1        1.2     1.3         2.1    2.2 2.3 3.1      3.2      3.3      4.1      4.2      4.3        5.1          5.2      5.3
subbeat  1.1.1  1.1.2                 2.1.1  2.1.2    3.1.1 3.1.2      3.3.1 3.3.2  4.1.1 4.1.2 4.2.1 4.2.1      5.1.1  5.1.2    5.2.1 5.2.2

**beamed**



**unbeamed**



**grouping** notes with measure bars and beams

- eases **readability** (player reads in a real-time context)

- highlight the **metric structure**
  hierarchy of strong / weak beats

# Common Western Music Notation

**metric structure**

bar  1  2  3  4  5

beat  1.1  1.2  1.3  2.1  2.2  2.3  3.1  3.2  3.3  4.1  4.2  4.3  5.1  5.2  5.3

subbeat  1.1.1  1.1.2  2.1.1  2.1.2  3.1.1  3.1.2  3.3.1  3.3.2  4.1.1  4.1.2  4.2.1  4.2.1  5.1.1  5.1.2  5.2.1  5.2.2

durations:
$$\frac{1}{2}\frac{1}{4}\frac{1}{4} \quad \frac{1}{16}\frac{1}{16}\frac{3}{4} \quad \frac{1}{16}\frac{1}{16}\frac{3}{4} \quad 0\ \frac{1}{2}\frac{1}{4}\frac{1}{4} \quad 2 \quad \frac{1}{2}\frac{1}{4}\frac{1}{4} \quad \frac{1}{16}\frac{1}{16}\frac{3}{4} \quad \frac{1}{2}\frac{1}{4}\frac{1}{4} \quad \frac{1}{2}\frac{1}{4}\frac{1}{4} \quad \frac{1}{2}\frac{1}{4}\frac{1}{4} \quad \frac{1}{16}\frac{1}{16}\frac{3}{4} \quad \frac{1}{2}\frac{1}{6}\frac{1}{6}\frac{1}{6} \quad \frac{1}{2}\frac{1}{2} \quad 0\ 1$$

Tree representation of the proportional rhythmic notation
with hierarchical encoding of durations: "*the* (duration) *data is in the structure*"
• the tree leaves contain the events
• the branching define durations, by partitioning of time intervals



single event
(note)

beamed division
arity = 2

$b_2$

$b_2$

$b_2$

corresponding
timeline

$0$     $\dfrac{1}{2}$   $\dfrac{3}{4}$   $1$

continuation
of the previous event

$b_2$

$b_2$    $-$

$b_2$

$b_2$

$0$ $\dfrac{1}{8}$ $\dfrac{1}{4}$   $\dfrac{1}{2}$    $1$

$b_2$

$b_2$    •

$b_2$

•$2$   •

1 grace-note (duration = 0)
and 1 note

$0$   $\dfrac{1}{4}$ $\dfrac{1}{2}$    $1$

defined by a Regular Tree Grammar:

- non-terminal symbols: $q, q_0, q_1, \ldots$

- terminal symbols (constants): • (1 note), $\bullet_2$ (1 grace-note + 1 note), — (continuation)

- production rules:

$q \rightarrow m_2(q_0, q) \mid m_0$

$q_0 \rightarrow u_3(q_1, q_1, q_1) \mid \bullet$      measure

$q_1 \rightarrow b_2(q_2', q_2) \mid \bullet \mid \bullet_2 \mid -$      beat = ♩

$q_2' \rightarrow b_2(q_3', q_3) \mid \bullet \mid \bullet_2 \mid -$    $q_2 \rightarrow b_2(q_3, q_3) \mid \bullet \mid -$    sub-beat = 8th-note = ♪

$q_3' \rightarrow \bullet \mid \bullet_2 \mid -$      $q_3 \rightarrow \bullet \mid -$      sub-sub-beat = 16nth note = ♬

derivations (lefmost)

$q_1 \rightarrow b_2(q_2', q_2) \rightarrow b_2(b_2(q_3', q_3), q_2) \rightarrow b_2(b_2(\bullet_2, q_3), q_2) \rightarrow b_2(b_2(\bullet_2, \bullet), q_2) \rightarrow b_2(b_2(\bullet_2, \bullet), \bullet)$



$q \rightarrow m_2(q_0, q) \rightarrow m_2(u_3(q_1, q_1, q_1), q) \rightarrow m_2(u_3(b_2(q_2', q_2), q_1, q_1), q) \rightarrow m_2(u_3(b_2(\bullet, q_2), q_1, q_1), q) \rightarrow \ldots$

**piano roll**
= sequence of timestamped input events

*structuring a linear representation
according to a language model*   =   **parsing**

tree-structured representation
of an output **music score**

conforming to a
prior language (expected notation)

2 nested extensions of parsing are needed
for the case music transcription:
- weighted extension
- symbolic weighted extension
  (*quantitative parsing*)

terminal symbols: $e_0, \ldots$ in a finite alphabet

input sequence

$$e_0 \, e_1 \qquad \ldots \qquad e_n$$

equality

yield
(sequence of leaves)

$$e_0 \, e_1 \qquad \ldots \qquad e_n$$

Parse Tree

parse-tree = representation of a
leftmost derivation of $e_0 \, e_1 \ldots e_n$ by a prior CF-grammar $\mathcal{G}$
with production rules: $q_0 \rightarrow q_1 \, q_2$ or $q_0 \rightarrow e$
(non-terminal symbols: $q_0, q_1, \ldots$ )

Decision problem: (membership)
does there exists a parse tree (leftmost derivation) of $\mathcal{G}$
that yields $e_0 \, e_1 \ldots e_n$ ?

Returning a parse tree of $\mathcal{G}$ that yields $e_0\,e_1\ldots e_n$

input sequence

$e_0\,e_1 \qquad \ldots \qquad e_n$

equality

yield
(sequence of leaves)

$e_0\,e_1 \qquad \ldots \qquad e_n$

Parse Tree

With an ambiguous prior CF-grammar $\mathcal{G}$ there might exists several parse trees (exponentially many).

in order to choose one (or some) parse trees, rank them according to their weight values, computed by Weighted Tree Grammar

**Weighted Regular Tree Grammar** $\mathcal{G}$:

- non-terminal symbols: $q, q_0, q_1, \ldots$
- terminal symbols (constants): $\bullet$ (1 note), $\bullet_2$ (1 grace-note + 1 note), $-$ (continuation)
- every production rule is assigned a weight value (*e.g.* cost to read):

$$q \xrightarrow{0} \mathsf{m}_2(q_0, q) \qquad q \xrightarrow{0} \mathsf{m}_0$$

$$q_0 \xrightarrow{0.1} \mathsf{u}_3(q_1, q_1, q_1) \qquad q_0 \xrightarrow{1} \bullet$$

measure

$$q_1 \xrightarrow{0.1} \mathsf{b}_2(q_2', q_2) \qquad q_1 \xrightarrow{1} \bullet \qquad q_1 \xrightarrow{1.9} \bullet_2 \qquad q_1 \xrightarrow{1} -$$

beat = ♩

$$q_2' \xrightarrow{0.1} \mathsf{b}_2(q_3', q_3) \qquad q_2' \xrightarrow{1} \bullet \qquad q_2' \xrightarrow{2.25} \bullet_2 \qquad q_2' \xrightarrow{1} -$$

sub-beat = 8th-note = ♪

$$q_2 \xrightarrow{0.1} \mathsf{b}_2(q_3, q_3) \qquad q_2 \xrightarrow{1} \bullet \qquad q_2 \xrightarrow{1} -$$

$$q_3' \xrightarrow{1} \bullet \qquad q_3' \xrightarrow{3.25} \bullet_2 \qquad q_3' \xrightarrow{1} - \qquad q_3 \xrightarrow{1} \bullet \qquad q_3 \xrightarrow{1} -$$

sub-sub-beat = 16th note = ♬

derivation (lefmost): $d : q_1 \xrightarrow{0.1} \mathsf{b}_2(q_2', q_2) \xrightarrow{0.1} \mathsf{b}_2(\mathsf{b}_2(q_3', q_3), q_2) \xrightarrow{3.25} \mathsf{b}_2(\mathsf{b}_2(\bullet_2, q_3), q_2) \xrightarrow{1} \mathsf{b}_2(\mathsf{b}_2(\bullet_2, \bullet), q_2) \xrightarrow{1} \mathsf{b}_2(\mathsf{b}_2(\bullet_2, \bullet), \bullet)$

cost of derivation: $\mathsf{weight}(d) = 0.1 + 0.1 + 3.25 + 1 + 1$

learning weight values from corpus statistics
Francesco Foscarin

In general, the weight values are taken in a commutative Semiring $\langle \mathbb{S}, \oplus, \mathbb{O}, \otimes, \mathbb{1} \rangle$

- $\oplus$ and $\otimes$ are associative and commutative, with neutral elements $\mathbb{O}$ and $\mathbb{1}$

- $\otimes$ distributes over $\oplus$ : $x \otimes (y \oplus z) = (x \otimes y) \oplus (x \otimes z)$

- $\mathbb{O}$ is absorbing for $\otimes$ : $\mathbb{O} \otimes x = \mathbb{O}$

| | domain | $\oplus$ | $\otimes$ | $\mathbb{O}$ | $\mathbb{1}$ |
|---|---|---|---|---|---|
| Boolean | $\{\bot, \top\}$ | $\vee$ | $\wedge$ | $\bot$ | $\top$ |
| Viterbi | $[0,1] \subset \mathbb{R}$ | max | $\times$ | $0$ | $1$ |
| Tropical min-plus | $\mathbb{R}_+ \cup \{+\infty\}$ | min | $+$ | $+\infty$ | $0$ |

Moreover, $\oplus$ is assumed to extend to infinite sums: there is an operation $\bigoplus\limits_{i \in I} x_i$ for all $I \subseteq \mathbb{N}$ such that:

*infinite sums extend finite sums*: $\forall j, k \in \mathbb{N}, j \neq k, \bigoplus\limits_{i \in \varnothing} x_i = \mathbb{O}, \bigoplus\limits_{i \in \{j\}} x_i = x_j, \bigoplus\limits_{i \in \{j,k\}} x_i = x_j \oplus x_k$

*associativity* and *commutativity*:

for all partition $(I_j)_{j \in J}$ of $I$, $\bigoplus\limits_{j \in J} \bigoplus\limits_{i \in I_j} x_i = \bigoplus\limits_{i \in I} x_i$

*distributivity* of products over infinite sums: for all $I \subseteq \mathbb{N}$, $\forall x, y \in \mathbb{S}$

$\bigoplus\limits_{i \in I} (x \otimes y_i) = x \otimes \bigoplus\limits_{i \in I} y_i$ and $\bigoplus\limits_{i \in I} (x_i \otimes y) = (\bigoplus\limits_{i \in I} x_i) \otimes y$

| | domain | $\oplus$ | $\otimes$ | $\mathbb{0}$ | $\mathbb{1}$ |
|---|---|---|---|---|---|
| Boolean | $\{\perp, \top\}$ | $\vee$ | $\wedge$ | $\perp$ | $\top$ |
| Viterbi | $[0,1] \subset \mathbb{R}$ | $\max$ | $\times$ | $0$ | $1$ |
| Tropical min-plus | $\mathbb{R}_+ \cup \{+\infty\}$ | $\min$ | $+$ | $+\infty$ | $0$ |

$\otimes$ is for composition of rule's weights in derivations and $\oplus$ is for optimal choice:

For a Weighted Regular Tree Grammar $\mathscr{G}$

$$\text{weight}_{\mathscr{G}}(d : q \xrightarrow{w_1} \ldots \xrightarrow{w_n} t) = \bigotimes_{i=1}^{n} w_i \quad \text{and} \quad \text{weight}_{\mathscr{G}}(q, t) = \bigoplus_{d:q \xrightarrow{+} t} \text{weight}_{\mathscr{G}}(d)$$

or recursively:

$$\text{weight}_{\mathscr{G}}(q, a(t_1, \ldots, t_n)) = \bigoplus_{q \xrightarrow{w} a(q_1, \ldots, q_n) \in \mathscr{G}} \left( w \otimes \bigotimes_{i=1}^{n} \text{weight}_{\mathscr{G}}(q_i, t_i) \right)$$

| | domain | $\oplus$ | $\otimes$ | $\mathbb{O}$ | $\mathbb{I}$ |
|---|---|---|---|---|---|
| Boolean | $\{\perp, \top\}$ | $\vee$ | $\wedge$ | $\perp$ | $\top$ |
| Viterbi | $[0,1] \subset \mathbb{R}$ | max | $\times$ | $0$ | $1$ |
| Tropical min-plus | $\mathbb{R}_+ \cup \{+\infty\}$ | min | $+$ | $+\infty$ | $0$ |

$\mathbb{S}$ is assumed :

- idempotent  $x \oplus x = x$

  that induces a partial ordering:  $x \leq_\oplus y$  iff $x \oplus y = x$

- total :  $\forall x, y \in \mathbb{S}$, either $x \oplus y = x$ or  $x \oplus y = y$    *i.e.* $\leq_\oplus$ is total

- bounded : $\mathbb{I} \oplus x = \mathbb{I}$,  or equivalently:  $\forall x, y \in \mathbb{S}, x \leq_\oplus x \otimes y$

  *i.e.* combining elements with $\otimes$ always increases their weight,
  see the *non-negative weights* condition for Dijkstra's shortest path algorithm
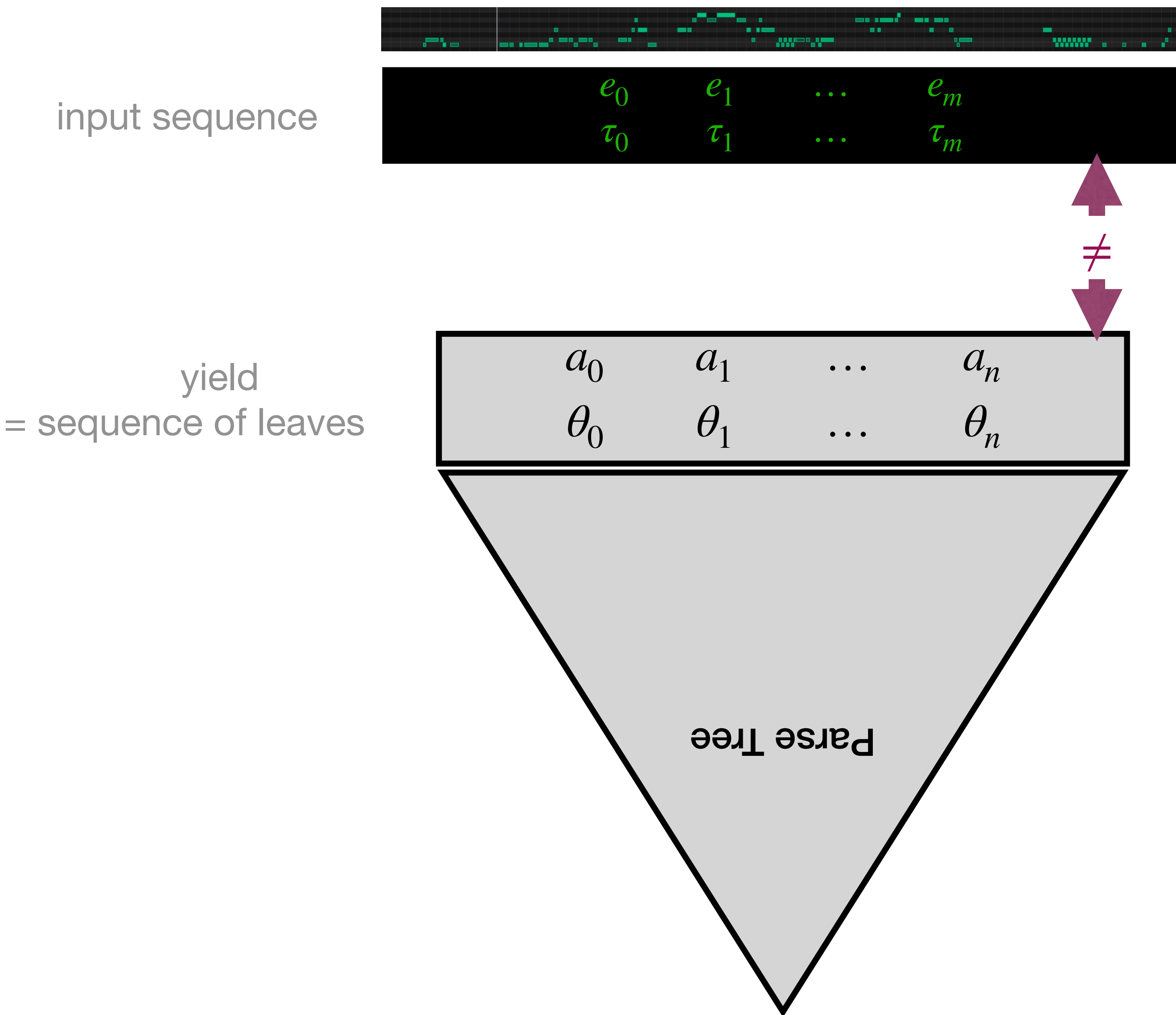
$k$-best parsing : enumeration of the $k$ best weighted trees *wrt* $\leq_\oplus$ for $\mathscr{G}$ and a non-terminal $q$, in PTIME,
user the above assumptions.

Similar to best path search in hyper-graphs (Dynamic Programming)
-   Viterbi algorithm in acyclic case
-   Knuth generalization of Dijkstra's algorithm in the general case

there is no 1-1 correspondance between input sequence and output leave sequence

input sequence

$$e_0 \quad e_1 \quad \ldots \quad e_m$$
$$\tau_0 \quad \tau_1 \quad \ldots \quad \tau_m$$

$\neq$

yield
= sequence of leaves

$$a_0 \quad a_1 \quad \ldots \quad a_n$$
$$\theta_0 \quad \theta_1 \quad \ldots \quad \theta_n$$

Parse Tree

we extend weighted parsing by ranking solutions with:

a measure of input / output fitness
= cost of IO alignement

$\otimes$

measure of cost-to-read
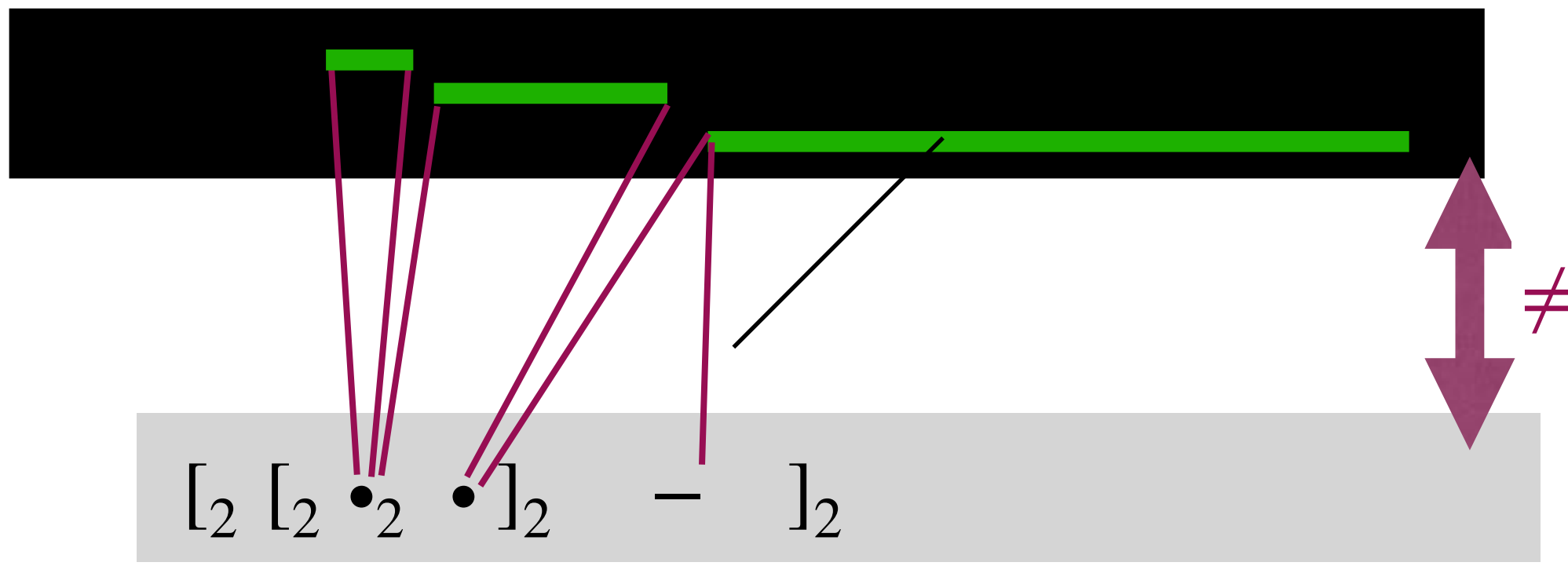weight value
computed by the Weighted Tree Grammar

| $E_{on}$ | $E_{off}$ | $D_{on}$ | $D_{off}$ | $C_{on}$ | $C_{off}$ |
|---|---|---|---|---|---|
| 0.11 | 0.19 | 0.22 | 0.48 | 0.53 | 1.08 |

input sequence

$\neq$

linearisation
of the output tree

$[_2 \quad [_2 \quad \bullet_2 \quad \bullet]_2 \quad - \quad ]_2$

cost of IO alignement
computed by a
Weighted word-to-word Transducer
(stateful definition of an edit-distance)

$b_2$

$b_2 \qquad -$

$\bullet_2 \qquad \bullet$

$0 \qquad \frac{1}{4} \qquad \frac{1}{2} \qquad 1$

input        output
symbol    symbol        desynchronise

$q_0 \xrightarrow{\langle E_{on}, \bullet_2 \rangle} q_1 \xrightarrow{\langle E_{off}, \varepsilon \rangle} q_2 \xrightarrow{\langle D_{on}, \varepsilon \rangle} q_3$

$\xrightarrow{\langle D_{off}, \bullet \rangle} q_4 \xrightarrow{\langle C_{on}, \varepsilon \rangle} q_5 \xrightarrow{\langle \varepsilon, - \rangle} q_5$

| $E_{on}$ | $E_{off}$ | $D_{on}$ | $D_{off}$ | $C_{on}$ | $C_{off}$ |
|----------|-----------|----------|-----------|----------|-----------|
| 0.11 | 0.19 | 0.22 | 0.48 | 0.53 | 1.08 |

grace-rests (eliminated): OFF and ON aligned to the same point

input sequence

$\neq$

linearisation of the output tree

$[_2\ [_2\ \bullet_2\ \bullet]_2\ \ -\ \ ]_2$

$b_2$

$b_2$ $-$

$\bullet_2$ $\bullet$

$\bullet_2$ $\bullet$

0    $\dfrac{1}{4}$   $\dfrac{1}{2}$    1

cost of IO alignement
computed by a
Weighted word-to-word Transducer
(stateful definition of an edit-distance)

input symbol    output symbol    desynchronise
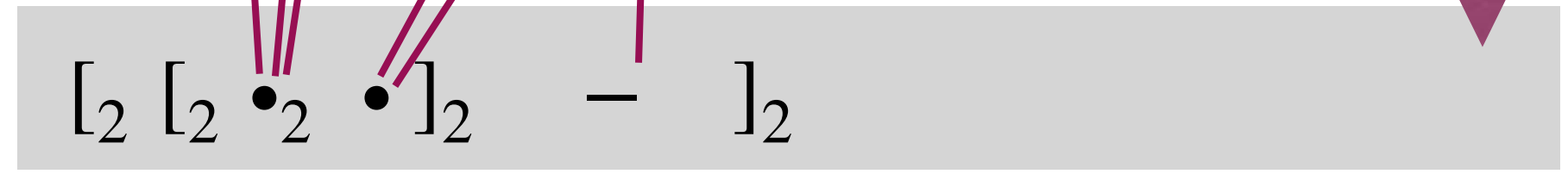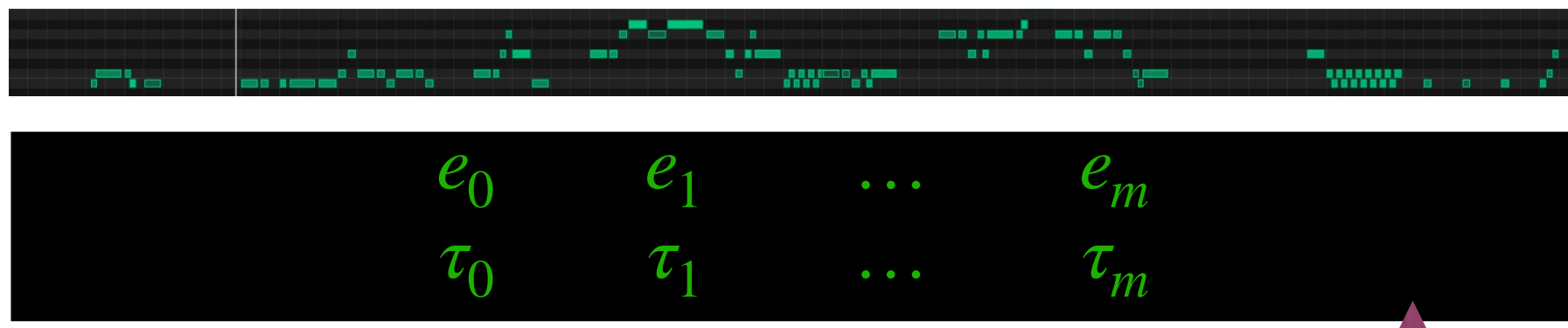
$q_0 \xrightarrow{\langle E_{on}, \bullet_2 \rangle} q_1 \xrightarrow{\langle E_{off}, \varepsilon \rangle} q_2 \xrightarrow{\langle D_{on}, \varepsilon \rangle} q_3$

$\xrightarrow{\langle D_{off}, \bullet \rangle} q_4 \xrightarrow{\langle C_{on}, \varepsilon \rangle} q_5 \xrightarrow{\langle \varepsilon, - \rangle} q_5$

in the context of music transcription, the symbols are timestamped → infinite alphabet $\Sigma_{inf}$
the weighted formalisms below must be able to read such symbols → symbolic extension

input sequence

$$e_0 \quad e_1 \quad \ldots \quad e_m$$
$$\tau_0 \quad \tau_1 \quad \ldots \quad \tau_m$$

measure of input / output fitness
= cost of IO alignement
computed by a
Weighted word-to-word Transducer

yield
= sequence of leaves

$$a_0 \quad a_1 \quad \ldots \quad a_n$$
$$\theta_0 \quad \theta_1 \quad \ldots \quad \theta_n$$

decorated with dates $\theta$
(computed with the durations
encoded in the tree structure)

$\otimes$

measure of cost-to-read
computed by the Weighted Tree Grammar

Parse Tree

Symbolic Weighted Language Models

SW-A: $\Sigma^*_{\text{inf}} \to \mathbb{S}$

$q \xrightarrow{\phi} q'$

$\phi : \Sigma_{\text{inf}} \to \mathbb{S}$

Weighted-A: $\Sigma^*_{\text{fin}} \to \mathbb{S}$

$q \xrightarrow{a,w} q'$

$a \in \Sigma_{\text{fin}}, w \in \mathbb{S}$

Droste, M., Kuich, W., Vogler
Handbook of WA, 2009

Symbolic-A: $\Sigma^*_{\text{inf}} \to \mathbb{B}\text{ool}$

$q \xrightarrow{\phi} q'$

$\phi : \Sigma_{\text{inf}} \to \mathbb{B}\text{ool}$

Veanes et al.
CAV 2017, CACM 2021

NFA: $\Sigma^*_{\text{fin}} \to \mathbb{B}\text{ool}$

$q \xrightarrow{a} q'$

$a \in \Sigma_{\text{fin}}$

for the transformation of the intermediate score representation



questions: rewrite strategies (*e.g.* IO or OI), conflicts…

- Piano transcription system (Kyoto U.)

  Non-local musical statistics as guides for audio-to-score piano transcription
  Kentaro Shibata, Eita Nakamura, Kazuyoshi Yoshii

  - deep-neural-network-based multipitch detection
    audio to unquantized MIDI

  - statistical-model-based (HMM) rhythm quantization
    unquantized MIDI to quantized MIDI

  - delegate to Muse Score + Voice separation algorithm for
    quantized MIDI to score

  - study of use of non-local statistics (pitch and rhythm)
    for the inference of global characteristics (metre, bar line positions…)

- Score Transformer (Yamaha) - piano transcription

  Score Transformer: Generating Musical Score from Note-level Representation
  Masahiro Suzuki

  Transformer model trained with popular songs (piano arrangements), KernScores (piano Sonata)
  MIDI to score (tokenization)

Implementation (FJ) of

- the above transcription by parsing framework

- the intermediate score model (w. Philippe Rigaux)

- other subtasks: pitch-spelling, key estimation, beat tracking…

https://gitlab.inria.fr/qparse/qparselib
https://qparse.gitlabpages.inria.fr

**qparse**: 75 Kloc C++

- command lines tools:
  `monoparse`, `drumparse`, grammar-learning, engraving (from quantified input)

- Python binding - Lydia Rodrigez-de la Nava
  scripts for automatic evaluation

- online port, real-time - Leyla Villaroel

**OpenMusic RQ lib**
Adrien Ycart 2016-17
IRCAM

https://forge.ircam.fr/p/omlibraries/downloads/
http://repmus.ircam.fr/cao/rq

CommonLisp/CLOS, 350 functions, 4900 lines of code

UI: Open Music object, input from chord-seq (notes, onset, dur) + segmentation marks,
      output to voice (OM rhythm trees)

# Score diff

by Francesco Foscarin
- identify the diff. between 2 XML music scores
- string/tree edit distance applied to a intermediate score representation

**Lamarque-Goudard** dataset (w. Francesco Foscarin, Teysir Baoueb)

- 283 monophonic extracts
  of classical repertoire
  inspired by a rhythm learning method

- ~ 20 measures per extract

- progressive difficulty
  cover a very large spectrum of rhythmic features

- score files (XML) and MIDI performances
  for evaluation and calibration of transcription tools

**Generation of artificial performances**
Madoka Goto, Masahiko Sakai (Nagoya U.), Satoshi Tojo (JAIST)

- construction of a GTTM tree

- segmetation accordingly

- performance generation by Director Musices (Anders Friberg)

monophonic : one note at a time
Good results for complex cases (ornaments, mixed tuplets, mixed note durations, silences…)
~ 100ms for the transcription of 1 score

Polonaise in D minor from Notebook for Anna Magdalena
Bach   BWV Anh II 128

original score

transcription of MIDI recording by qparse

Polonaise in D minor from Notebook for Anna Magdalena Bach   BWV Anh II 128

original score

transcription of MIDI recording by Finale

Beethoven, Trio for violin, cello
and piano, op.70 n.2 (2d mov)

original score



transcription
of MIDI recording
by qparse

Beethoven, Trio for violin, cello
and piano, op.70 n.2 (2d mov)

original score



transcription
of MIDI recording
by Finale

options:
- mixed rhythms,
- tuplets
- smallest note = 32nd
The time signature
and the tempo are given.

**FiloBass** by John-Xavier Riley (QMUL, C4DM)
project "*Dig That Lick*"

- jazz bass lines, acc. of saxophone

- 48 tracks,
  24 recorded hours of melodies and improvisations

- qparse as backend of an audio-to-MIDI
  transcription procedure

- prior beat (measure) tracking

**Groove MIDI** Dataset

- by Google Magenta

- 13.6 hours, 1150 MIDI files, ~ 22000 measures recorded by professional drummers on a electronic drum kit

- audio (wav) files synthesized from (and aligned to) MIDI files for evaluation of audio-to-MIDI drum transcription

- no score files!

**Scoring the GMD** with qparse

Martin Digard (INALCO)

- all score files (XML) produced from the MIDI files with the same generic tree grammar (4/4 measure)

- polyphonic case-study, simpler than piano

- specific drumming constraints (hands ≤ 2, feet ≤ 2)

- processing errors from MIDI sensors

- Dataset **ASAP** - Francesco Foscarin, Andrew Mc Leod
  MIDI and audio recording from Yamaha piano competition
  + XML scores
  + alignments
  + beat tracking annotations

- **Voice separation** - Lydia Rodrigez-de la Nava, evaluation Augustin Bouquillard
  and for piano guitar transcription.
  integration in transcription:

  - before parsing, or

  - after parsing (on intermediate model), or

  - joint with parsing.

**MIDI-to-Score Automated Music Transcription approach**

- quantitative parsing technique
  based on Symbolic Weighted formal language formalisms
  Tree Automata and word-to-word Transducers

- with prior quantitative language of notation *style*
  and prior IO measure

- (abstract) hierarchical score model
  as intermediate representation for score generation

- can handle complex notation cases:
  ornaments, mixed tuplets, mixed note durations, silences…

- efficient

- case studies: Monophonic, Drums

- ongoing work on Polyphonic case studies: guitar, piano

# MERCI!
# THANK YOU!