

Tempo curves estimation, generation and analysis, L3 intership (CNAM / INRIA)

27/05/24 - 02/08/24

Sylvain Meunier (intern)
sylvain.meunier@ens-rennes.fr

Florent Jacquemard (supervisor)
florent.jacquemard@inria.fr

Abstract— Tempo estimation consists in detecting the speed at which a musician plays, or more generally at which a piece of music is played or heard. Since tempo may not be constant at the scale of a piece, even locally, we need some kind of reference to compare to in order to define said speed. Indeed, a note-wise speed would not match the intuitive notion of tempo, based on a regular *pulse*. Such a reference can be found in Western symbolic notations of music, called either *music score* or *sheet music*, that allows for a definition of tempo as symbolic speed. We present here some results regarding the generation and analysis of local tempo curves of musical performances, involving methods that need to be given some symbolic information, and methods that generate them on the fly.

Index terms—Music Information Retrieval, tempo estimation, quantification, musical formalism

I. INTRODUCTION

The [Music Information Retrieval \(MIR\)](#) community focuses on three representations of musical information. The first one is raw audio, either recorded or generated, encoded using WAVE or MP3 format. The computation is based on a physical understanding of signals, using audio frames and spectrum, and represents the most common and accessible type of data. The second is a more musically-informed format, representing notes with both pitch (i.e., the note that the listener hear) and duration, encoded within a MIDI file. [Such a file can be displayed as a piano roll, i.e., a graph whose x-axis is time and y-axis is pitch (hence, the y-axis is discrete)]. The last way to encode musical information is the computed counterpart of sheet music. A sheet music is a way to write down a musical score, that is usually computed as a MusicXML file, mainly for display and analysis purposes. It comes with a *symbolic* and abstract notation for time, that only describes the length of events in relation to a specific abstract unit, called a *beat*, and the pitch of each event. These data are actually ordered from the most to the least common and accessible on a large scale.

To actually play a sheet music, one needs a given [tempo](#), usually indicated as the amount of beat per minute (BPM).

Therefore, the notion of tempo allows to translate symbolic notation expressed in [MTU](#) to real time events expressed in [RTU](#). We will discuss later on a formal definition of tempo. However, tempo itself is insufficient to describe an actual performance of a sheet music, i.e., the sequence of real time events. Indeed, S. D. Peter, C. E. Cancino-Chacón, E. Karystinaios, and G. Widmer [1] present four parameters, among which tempo and [articulation](#) appear the most salient in contrast with [velocity](#) and timing. The latter represents the delay between the theoretical real time onset according to the current tempo, and the actual onset heard in the performance. Even though such a delay is inevitable for neurological and biological reasons, those timings are usually overemphasized and understood as part of the musical expressivity of the performance.

Unlike one may imagine, a performance of a given sheet music does need to contains some errors in order to sound “musical”. Otherwise, the result would sound too unhumanize, automatic, i.e., unalive and thus unmusical.

Even though the community studies the four parameters, the hierarchy exposed by [1] embodies quite well their relative priority within litterature. O. F. B. Katerina Kosta Rafael Ramírez and E. Chew [2] present results pointing that, although velocities don’t help to meaningfully estimate tempo, the latter allows to marginally upgrade velocity-related predictions. Actually, velocity appears to be more of a score parameter rather than a performance one : automatic learning methods trained on performances of a single piece showed much better results when asked to predict velocities employed by another performer on the same piece than when trained on other performances of the same performer.

In this study, we shall focus mainly on tempo estimation for a given performance recorded as a MIDI file, on both a local and global level. Section II presents an overview of the domain and related problems. We then present the following contributions :

- a justified proposition for a formal definition of tempo based on [3], [4] and [5] ([III.A](#)) ; and some immediate consequences ([III.B](#))

- a revision of [6] and [7] to fit a score-based approach (III.C)
- a general theoretical framework for scoreless tempo estimation with application to [8], [9] (IV.B)
- an extension of [9], to fit tempo estimation, and some new theoretical results (IV.D, IV.E and Appendix C)
- A method for data augmentation, and results based on [10] and [11] (Section V.A)

This report and associated algorithms can be found on [12].

II. STATE OF ART

Tempo and related works actually hold a prominent place in litterature. Direct tempo estimation was first computed based on probabilistic models (C. Raphael [3], E. Nakamura, N. Ono, S. Sagayama, and K. Watanabe [13], [14]), and physical / neurological models (E. W. Large and M. R. Jones [6], H.-H. Schulze, A. Cordes, and D. Vorberg [7]) ; before the community tried neural network models [2] and hybrids approaches (K. Shibata, E. Nakamura, and K. Yoshii [15]). As the majority of previous examples, we shall focus here on mathematically and/or musically explainable methods.

Since tempo needs a symbolic representation to be meaningful, one can consider *transcription* as a tempo-related work. We will keep this discussion for Section IV.D and Section VI.

However, note-alignment, that is matching each note of a performance with those indicated by a given score is a very useful preprocessing technique, especially for direct tempo estimation and further analysis, such as [4], [16], [5]. Two main methods are to be found in litterature : a dynamic programming algorithm, equivalent to finding a shortest path (M. Müller [17]), that can works on raw audio (WAVE files) ; and a Hidden Markov Model (E. Nakamura, K. Yoshii, and H. Katayose [18]) that needs more formatted data, such as MIDI files.

III. SCORE-BASED APPROACHES

A. Formal considerations

Since we chose to focus on MIDI files, we will represent a performance as a strictly increasing sequence of events $(t_n)_{n \in \mathbb{N}}$, each element of whose indicates the onset of the corresponding performance event. Such a definition is very close to an actual MIDI representation.

For practical considerations, we will stack together all events whose distance in time is smaller than $\varepsilon = 20$ ms. This order of magnitude, calculated by E. Nakamura, T. Nakamura, Y. Saito, N. Ono, and S. Sagayama [14] represents the limits of human ability to tell two rhythmic events apart,

and is widely used within the field [15], [4], [16], [5], [18], [9]–[11], [8].

Likewise, a sheet music will be represented as a strictly increasing sequence of events $(b_n)_{n \in \mathbb{N}}$. In both of those definition, the terms of the sequence do not indicate the nature of the event (*chord*, single note, *rest*...). Moreover, in terms of units, (t_n) corresponds to real onset, thus expressed in seconds, whereas (b_n) corresponds to theoretical or symbolic onsets, expressed in beats.

With those definitions, let us formally define tempo $T(t)$ so that, for all $n \in \mathbb{N}$, $\int_{t_0}^{t_n} T(t) dt = b_n - b_0$.

Appendix A shows that this definition is equivalent to : $\forall n \in \mathbb{N}$, $\int_{t_n}^{t_{n+1}} T(t) dt = b_{n+1} - b_n$. However, tempo is only tangible (or observable) between two events *a priori*. We will then define the canonical tempo $T^*(t)$ so that :

$$\forall x \in \mathbb{R}^+, \forall n \in \mathbb{N}, x \in [t_n, t_{n+1}[\Rightarrow T^*(x) = \frac{b_{n+1} - b_n}{t_{n+1} - t_n}.$$

The reader can verify that this function is a formal tempo according to the previous definition. From now on, we will consider the convention : $t_0 = 0$ (s) et $b_0 = 0$ (beat).

Even though there is a general consensus in the field as for the interest and informal definition of tempo, several formal definitions coexist within litterature : K. Shibata, E. Nakamura, and K. Yoshii [15] and E. Nakamura, N. Ono, S. Sagayama, and K. Watanabe [13] take $\frac{1}{T}$ as definition ; C. Raphael [3], [4] et P. Hu and G. Widmer [5] choose similar definitions than the one given here (approximated at the scale of a *measure* or a section for instance).

T^* has the advantage to coincide with the tempo actually indicated on traditional sheet music (and therefore on MusicXML format), hence allowing a simpler and more direct interpretation of results.

B. Naive use of formalism

As said in introduction, the more formatted data, the less accessible it is ; and the field contains only a few datasets containing both sheet music and corresponding audio, more or less anotated with various labels [4], [16], [5], [10], [11].

In our study, we chose to rely on the (n)-ASAP dataset [11] that presents a vast amount of performances, with over 1000 different pieces of classical music, all note-aligned with the corresponding score. From there, we can easily compute our definition of tempo. Figure 1 presents the results for a specific piece of the (n)-ASAP dataset with a logarithmic y-scale, that contains a few brutal tempo change, whilst maintaining a rather stable tempo value in-between.

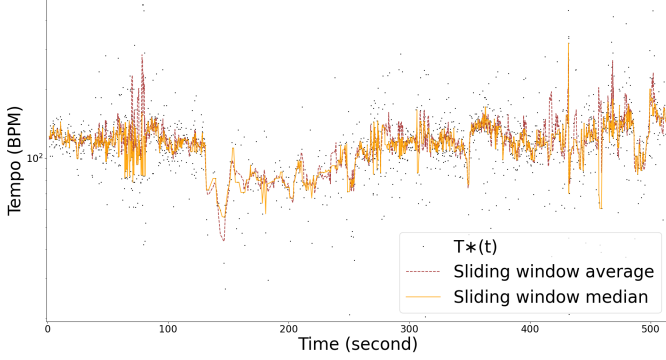


Figure 1: Tempo curve for a performance of Islamey, Op.18, M. Balakirev, with naive algorithm

In this graph, one can notice how T^* (plotted as little dots) appears noisy over time; even though allowing to distinguish a tempo change at $t_1 = 130$ s and $t_2 = 270$ s. Both the sliding window average (dotted line) and median (full line) of T^* seem unstable, presenting undesirable peaks, whereas the “feeled” tempo is quite constant for the listener, although the median line is a bit more stable than the average line, as expected. There are two explanation for those results. First, fast events are harder to play exactly on time, and the very definition being a ratio with a small theoretical value as the denominator explains the deviation and absurd immediate tempo plotted. In fact, we can read that about 10 points are plotted over 400 BPM (keep in mind that usual tempo are in the range 40 - 250 BPM). Second, the notion of timing and tempo are mixed together in this computation, hence giving results that do not match the listener feeling of a stable tempo. Actually, timing can be seen as little modifications to the “official” score, and using the resulting score would allow for curves that fit better the listener feeling, though needing an actual transcription of the performance first.

C. Physical models

Among the tasks needing tempo estimation, the problem of real time estimation to allow a dedicated machine to play an accompagnement by following at least one real musician has been tackled by various approaches in litterature. C. Raphael [3] started with a probabilistic model, but those methods have found themselves replaced by a more physical understanding of tempo *via* the notion of internal pulse, as explained by E. W. Large and M. R. Jones [6]. In fact, their method has recently been developped to a commercial form¹, based on an a previous adaption of the latters by A. Cont [19].

The approach developped by E. W. Large and M. R. Jones [6] consider a simplified neurological model, where listening is a fundamentally active process, implying a synchro-

nization between external events (those of the performance) and an internal oscillator, whose complexity depends of hypothesis on the shape of the first ones. The model consists of two equations for the internal parameters:

$$\Phi_{n+1} = \left[\Phi_n + \frac{t_{n+1} - t_n}{p_n} - \eta_\Phi F(\Phi_n) \right] \bmod_{[-0.5, 0.5]} 1 \quad (1)$$

$$p_{n+1} = p_n (1 + \eta_p F(\Phi_n)) \quad (2)$$

Here, (Φ_n) corresponds to the phase, or rather the phase shift between the oscillator and the external events, and (p_n) embodies its period. Finally, η_p and η_Φ are both constant parameters. This initial model is then modified to consider a notion of attending *via* the κ parameter, whose value change over time according to other equations. The new model contains the same formulas, with the following definition for F

$$F : \Phi, \kappa \rightarrow \frac{\exp(\kappa \cos(2\pi\Phi)) \sin(2\pi\Phi)}{\exp(\kappa)} \cdot \frac{1}{2\pi}.$$

Even though this model shows pretty good results, has been validated through some experiments in [6], and is still used in the previously presented version (E. W. Large *et al.* [20]), a theoretical study of the system behavior remains quite complex, even in simplified theoretical cases [7], notably because of the function F expression.

In order to simplify the previous model, H.-H. Schulze, A. Cordes, and D. Vorberg [7] present *TimeKeeper*, that can be seen as a linearization of the previous approach, valid in the theoretical framework of a metronome presenting small tempo variations. In fact, there is a strong analogy between the two models, that are almost equivalent under specific circumstances, as shown by J. D. Loehr, E. W. Large, and C. Palmer [21]. Here, we used the derandomised version presented in [21], where $M_i = 0$ and $T_i = \tau_i$ for all $i \in \mathbb{N}$.

None of those models have an inherent comprehension of musical score information, since they both rely on a rather stable metronome. In the version displayed hereafter, they were modified to consider score information, in order to create a more stable and precise value of tempo than the naive approach previously presented. Those modifications are detailed in the following paragraph (OR IN APPENDIX ?), and were made in order to keep consistency with the original models in their initial theoretical framework of validity. Let $\min_{\text{abs}} : a, b \mapsto \begin{cases} a & \text{if } |a| < |b| \\ b & \text{otherwise} \end{cases}$

We first modify the E. W. Large and M. R. Jones [6] equations accordingly :

$$\Phi_{n+1} = \Phi_n + \frac{t_{n+1} - t_n}{p_n} - \eta_\Phi F(\Psi_n, \kappa_n) \quad (3)$$

$$p_{n+1} = p_n (1 + \eta_p F(\Psi_n, \kappa_n)) \quad (4)$$

$$\Psi_n = -\min_{\text{abs}}(k + b_n - \Phi_n, k + 1 + b_n - \Phi_n)$$

$$k = \lfloor \Phi_n - b_n \rfloor$$

¹<https://metronautapp.com/>

Using the analogy presented in [21], we then obtain the following equations for *TimeKeeper* :

$$A_{i+1} = K_i(1 - \alpha) + \tau_i - (t_{i+1} - t_i) \quad (5)$$

$$\tau_{i+1} = \tau_i - \beta * (K_i \bmod_{[-0.5, 0.5]} 1) \quad (6)$$

$$K_i = -\min_{\text{abs}}(k\tau + b_i - A_i, (k+1)\tau + b_i - A_i)$$

$$k = \left\lfloor \frac{A_i - b_i}{\tau_i} \right\rfloor$$

In both of those extensions, the \min function is used in order to represent a choice between two corrections. The first argument can be interpreted as a correction with respect to the most recent passed beat time occurring exactly on a actual beat, i.e., $a_1 = \max_{n \in \mathbb{N} : b_n \leq B_i} [b_n]$ where B_i represents the internal value at time i acting as a beat unit (Φ_i for E. W. Large and M. R. Jones [6] and A_i for H.-H. Schulze, A. Cordes, and D. Vorberg [7]). The second argument embodies the correction according to the following beat, $a_2 = a_1 + 1$.

One can notice that the phase is actually always used modulo 1 in E. W. Large and M. R. Jones [6], since it appears only multiplied by 2π in either cos or sin functions. Using this remark, one can verify that, in the initial presentation of the model with a metronome, i.e., $\forall n \in \mathbb{N}, b_n = 0 \bmod 1$, the extension proposed here is equivalent to the original approach, i.e., (1), (2) \Leftrightarrow (3), (4), hence justifying the designation “extension”.

Figure 2 displays the results of those two models, in regards with the canonical, or immediate tempo. One can notice that E. W. Large and M. R. Jones [6] model is less stable than *TimeKeeper*, although faster to converge.

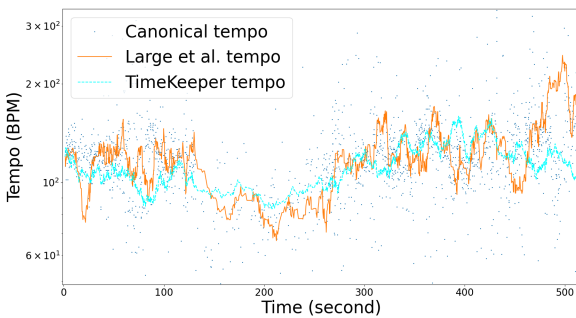


Figure 2: Tempo curve for a performance of Islamey, Op.18, M. Balakirev, according to various models

Figure 3 (below) exposes the visible difference in tempo initialization of the two models, starting both here with the initial tempo of 70 BPM ($\text{♩} = 70$, i.e., the *beat* unit here is a quarter note). *TimeKeeper* does not manage to converge to any significant tempo. Such a behavior was to be expected, considering the theoretical framework for *TimeKeeper*, that is

small tempo variation, and correct initialization. However, Large et al model manages to converge to a meaningful result. In fact, in the range 9 to 70 seconds, the estimated tempo according to Large is exactly half of the actual tempo hinted by the blue dots (canonical tempo).

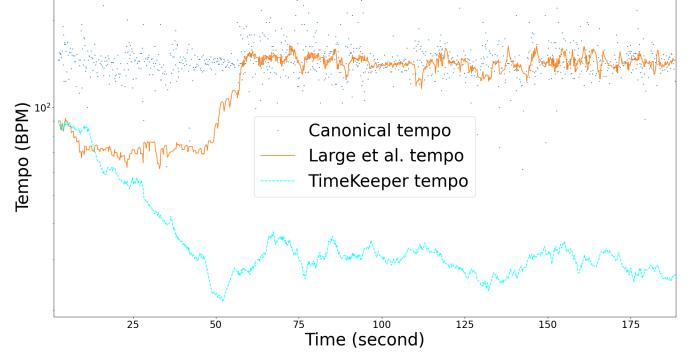


Figure 3: Tempo curve for a performance of Piano Sonata No. 11 in A Major (K. 331: III), W.A Mozart, according to the previous models with irrelevant initialization

IV. SCORELESS APPROACHES

A. Motivations

There are three main issues with the previous models, apart from the necessary knowledge of the sheet music, that are : salient sensibility to tempo initialization (cf. Figure 3), unstability that requires some time to (possibly) converge (cf. Figure 2 and Figure 3), and difficulty to accurately estimate relevant values of the constant internal parameters. According to our implementation, E. W. Large and M. R. Jones [6] is a particularly chaotic model regarding the latter.

We will present here two models focusing on tackling mainly the first two issues previously presented. Those rely on a specific musical property of division : in symbolic notations of music, every single event can be comprehended as a multiple of a certain unit called a *tatum*, usually expressed in beat unit. Therefore, the real events of a performance, or rather their duration, can be interpreted as multiple of this tatum. However, considering a non-constant tempo, the real value (i.e., real duration in seconds) of this tatum may evolve through time, whereas the symbolic value remains constant anyway. Actually, detecting the tatum is equivalent to transcribe the performance to sheet music, which a rather more complicated task than tempo estimation. For instance, there are several ways to write down sheet musics that are undistinguishable when performed. We call this ambiguity *tempo octaves* (cf Appendix A).

We considered here two quantization methods extracted from literature : D. Murphy [8] and G. Romero-García, C. Guichaoua, and E. Chew [9]. Both of these papers present a theory of approximate division, that is a way to find a ratio-

nal number, interpreted as the ratio of two symbolic events expressed in arbitrary unit (for instance in tatum) from the real events durations in seconds. Such a link is equivalent to defining a tempo with the formalism presented in III. Although [8] provides an algorithm to find candidate rationals, they do not include a way to compare those candidates, thus leaving no choice but an exhaustive approach (top-down in the paper). With another formalism, G. Romero-García, C. Guichaoua, and E. Chew [9] define a graph, restrained only to consistent values of tatum. We choose to adapt the latter, although their presentation is clearly user-oriented rather than automatic, since the introduced graph allows to mathematically (and therefore automatically) define a “good” and a “best” choice among all possible found tatum values.

B. Introduction of an estimator based approach

Given a sequence $(u_n)_{n \in \mathbb{N}}$, we now introduce the notation $(\Delta u_n)_{n \in \mathbb{N}} := (u_{n+1} - u_n)_{n \in \mathbb{N}}$ for the next sections.

The reason why the previous models have to converge is because they both try to find an exact value of tempo, and therefore sudden and huge tempo changes will lead to a uncertain period for the resulting tempo estimation, that is in this way the exact same problem as tempo initialization. When doing tempo estimation, we are in fact much more interested in a local tempo variation, relative to the previous estimation, rather than an absolute value, especially on a local time scale (where we can often assume tempo to be almost constant). Using the formalism presented in III, we first present the following result since $T_n^* > 0$:

$$T_{n+1}^* = T_n^* \frac{T_{n+1}}{T_n} = T_n^* \frac{\Delta t_n}{\Delta t_{n+1}} \frac{\Delta b_{n+1}}{\Delta b_n}$$

Let T_n be an estimation of T_n^* by a certain given model and $\alpha_n = \frac{T_n}{T_n^*}$. We obtain :

$$\alpha_n T_{n+1}^* = \underbrace{\alpha_n T_n^*}_{T_n} \frac{\Delta t_n}{\Delta t_{n+1}} \times \frac{\Delta b_{n+1}}{\Delta b_n}$$

In the above formula, the only value to actually estimate is therefore $\frac{b_{n+2}-b_{n+1}}{b_{n+1}-b_n}$, where both the numerator and denominator are translation permissive (i.e., we can afford a locally constant shift in both of our estimation), hence the resulting value is invariant by translation, or constant multiplication of our estimation. Furthermore, the value to be estimated only deals with symbolic units, meaning that we can use musical properties to find a consistent result. As a result, we can obtain a tempo estimation with the same multiplicative shift as the previous estimation T_n , thus, by using the formula recursively, we obtain a model that can track tempo variations over time without any need of convergence, and that is robust to tempo initialization, while using only local methods (in other words, the resulting model is [online](#)). As noticed in the beginning of this section, symbolic value have some bindings that actually help to determine their values.

The point of this approach is to keep a constant factor between (T_n) and (T_n^*) in order to prevent the need for any convergence time. We will now define $T_{n+1} = \alpha_n T_{n+1}^*$.

We then find : $\frac{\Delta b_{n+1}}{\Delta b_n} = \frac{T_{n+1} \Delta t_{n+1}}{T_n^* \Delta t_n} = \frac{\frac{1}{\alpha_n} T_{n+1}}{\frac{1}{\alpha_n} T_n} \times \frac{\Delta t_{n+1}}{\Delta t_n}$,

hence $\frac{\Delta b_{n+1}}{\Delta b_n} = \frac{T_{n+1}}{T_n} \times \frac{\Delta t_{n+1}}{\Delta t_n}$.

Let us then write the actual formula of the model :

$$\frac{T_{n+1}}{T_n} = \frac{T_{n+1}^*}{T_n^*} = \frac{\Delta t_n}{\Delta t_{n+1}} E \left(\underbrace{\frac{T_{n+1}}{T_n} \times \frac{\Delta t_{n+1}}{\Delta t_n}}_{\frac{\Delta b_{n+1}}{\Delta b_n}} \right) \quad (7)$$

where E is a function-like object (closer to a *object* in computer science than an actual mathematical function), designated by *estimator*. This function is supposed to act, on a theoretical ground, as an oracle that returns the correct value of the symbolic $\frac{\Delta b_{n+1}}{\Delta b_n}$ from the given real values indicated in (7), therefore supposed to rarely match the theoretical values.

Given an estimator E , the tempo value defined as T_{n+1} , computed from both T_n and local data, is obtained *via* the following equation, where x represents $\frac{T_{n+1}}{T_n}$ in (7) :

$$T_{n+1} = T_n \underset{x \in [\frac{\sqrt{2}}{2} T_n, \sqrt{2} T_n]}{\operatorname{argmin}} d \left(x, \frac{\Delta t_n}{\Delta t_{n+1}} E \left(x \frac{\Delta t_{n+1}}{\Delta t_n} \right) \right) \quad (8)$$

where $d : a, b \mapsto k_* |\log(\frac{a}{b})|$, $k_* \in \mathbb{R}_+^*$, is a logarithmic distance, chosen since an absolute distance would have favor small values by triangle inequality in the following process.

In the implementation presented here, the estimator role is more to quantify the ratio in order to output a musically relevant value. In our test, we limited these quantifications to accept only regular division (i.e., powers of 2). Furthermore, the numerical resolution for the previous equation was done by a logarithmically evenly spaced search and favor x values closer to 1 (i.e., T_{n+1} closer to T_n) in case of distance equality.

Such a research allows for a musically explainable result : the current estimation is the nearest most probable tempo, and both halving and doubling the previous tempo is considered as improbable, and as further going from the initial tempo. Appendix B gives further explanation about (8).

C. Study of the model

Since this approach fundamentally search to estimate tempo variation rather than actual values, it is not easy to visualize the relevance of the result by naive means. We choose here to define $\left(\alpha_n := \frac{T_n}{T_n^*} \right)_{n \in [1, N]}$ and $\left(\tilde{\alpha}_n := \exp \left(\ln \left(\frac{T_n}{T_n^*} \right) - \left\lfloor \log_2 \left(\frac{T_n}{T_n^*} \right) \right\rfloor \ln(2) \right) \right)_{n \in [1, N]}$. $(\tilde{\alpha}_n)_{n \in [1, N]}$ is then called the *normalized* sequence of ratio, where each value is uniquely determined within the range $[1, 2[$. Such a choice allows for merging together the tempo octaves, as explained in Appendix A. In this representation,

we actually have $\tilde{2} = \tilde{1}$, and adding $\tilde{1}$ is equivalent to multiply the initial value by 2. We will now define a *spectrum* $S = (\tilde{\alpha}_n)_{n \in \llbracket 1, N \rrbracket}$. We then call $|S|$ the value $N \in \mathbb{N}$ and \mathcal{C} the range $[1, 2[$ seen as a circle according to the following application : $x \mapsto (\cos(2\pi x), \sin(2\pi x))$, so that $c(1) = c(2)$

We then define the *measure* of a spectrum S , that embodies a standard deviation on \mathcal{C} , so that :

$$m(S, \Delta) = \max_{\tilde{d} \in \mathcal{C}} \frac{|\{n \in \llbracket 1, |S| \rrbracket : d(\tilde{\alpha}_n, \tilde{d}) \leq \Delta\}|}{|S|} \quad (9)$$

Where $\Delta \in [0, \frac{1}{2}]$ embodies the measure accuracy, and d is still a logarithmic distance, slightly modified on \mathcal{C} to be consistent with $d(\tilde{1}, \tilde{2}) = 0$. Actually, it can be shown that on \mathcal{C} , $d : \tilde{a}, \tilde{b} \mapsto \min(|\tilde{a} - \tilde{b}|, 1 - |\tilde{a} - \tilde{b}|)$.

The reader can verify that this measure is invariant with respect to spectrum rotation by any $\lambda \in R_+^*$ (i.e., $m(S, \Delta) = m(\left(\tilde{\lambda}\tilde{\alpha}_n\right)_{n \in \llbracket 1, N \rrbracket}, \Delta)$), and does not depend on the normalisation interval (here $[1, 2[$, but actually $[\lambda, 2\lambda[$ would work just as well, with a different expression of d on \mathcal{C}). Finally, $0 \leq m(S, \Delta) \leq 1$, and $m(S, \Delta) = 0 \Leftrightarrow |S| = 0$, $m(S, \Delta) = 1$ iff S only contains values within a 2Δ range.

This measure allows to quantify the quality of this model, without considering tempo octaves, or equivalently to quantify the quality of the estimator. The C++ implementation of the measure used to obtain the following figures if available on [12].

résultats sur (n-)Asap selon les périodes, les compositeurs, etc..., éventuellement en annexe ?

D. Towards a quantified approach

The previous model supposes the existence of an oracle, more or less correct in its predictions, that is actually a (partial) transcriber. In this section, we will focus on this transcribing part, by extending G. Romero-García, C. Guichaoua, and E. Chew [9] model with the previous formalism. In fact, we extend the previous approach by considering the estimator as our central model and then extracting tempo values rather than the opposite.

Let $n \in \mathbb{N}^*$ and $D \subset (R^+)^n$ be a set of some durations of real time events. [9] defined the continuous function ε_D as :

$$\varepsilon_D : a \mapsto \max_{d \in D} \min_{m \in \mathbb{Z}} |d - ma| \quad (10)$$

This function is called the *transcription error*, and can be interpreted as maximum error (in real time unit) between all real events $d \in D$ and theoretical real duration of ma , where m is a symbolic notation expressed in arbitrary symbolic unit, and a a real time value corresponding to a *tatum* at a given tempo. We prove in Appendix C that the set of all local maxima of ε_D , except those that also are minima, is :

$$\begin{aligned} M_D &= \left\{ \frac{d}{k + \frac{1}{2}}, d \in D, k \in \mathbb{N} \right\} \\ &= \bigcup_{d \in D} \left\{ \frac{d}{k + \frac{1}{2}}, k \in \mathbb{N} \right\} \end{aligned} \quad (11)$$

In fact, each of these local maxima corresponds to a change of the m giving the minimum in the expression of ε_D , hence the following result : in-between two such successive local maxima, the quantification remains the same, i.e.,

Proposition IV.D.1 Let m_1, m_2 be two successive local maxima of ε_D , $a_1 \in]m_1, m_2[$, $a_2 \in [m_1, m_2]$, $d \in D$ and $m \in \mathbb{Z}$. Then $m \in \operatorname{argmin}_{k \in \mathbb{Z}} |d - ka_1| \Rightarrow m \in \operatorname{argmin}_{k \in \mathbb{Z}} |d - ka_2|$.

With this property, we can then choose to consider only local minima of ε_D as in [9], since there is exactly one local minima in-between two such successive local maxima, and choosing any other value in this range would result in the exact same transcription, with a higher error by definition of a local maxima (that is global on the considered interval). The correctness of the following algorithm to find all local minima within a given interval is proven in Appendix C.

FindLocalMinima($D \neq \emptyset$, start, end) :

```

1   $M \leftarrow \left\{ \frac{d}{k + \frac{1}{2}}, d \in D, k \in \left[ \left\lfloor \frac{d}{\text{end}} \right\rfloor, \left\lfloor \frac{d}{\text{start}} \right\rfloor \right] \right\}$ 
2   $P \leftarrow \left\{ \frac{d_1 + d_2}{k}, (d_1, d_2) \in D^2, k \in \left[ \left\lfloor \frac{d_1 + d_2}{\text{end}} \right\rfloor, \left\lfloor \frac{d_1 + d_2}{\text{start}} \right\rfloor \right] \right\}$ 
3  localMinima  $\leftarrow \emptyset$ 
4  for  $(m_1, m_2) \in M^2$  two successive maxima :
5    | in_range  $\leftarrow \{p \in P : m_1 \leq p < m_2\}$ 
6    | localMinima  $\leftarrow \text{localMinima} \cup \{\min(\text{in\_range})\}$ 
7  return localMinima
```

Algorithm 1: Returns all local minima within [start, end]

G. Romero-García, C. Guichaoua, and E. Chew [9] then defined $G = (V, E)$ a graph whose vertices are the local minima of ε_D with D a sliding window, or *frame*, on a given performance, and whose edges are so that they can guarantee a *consistency property*, explained hereafter.

The *consistency property* for two tatums a_1, a_2 specifies that, if F_\cap is the set of all values in common between two successive frame, for all $d \in F_\cap$, d is quantified the same way according to the tatum a_1 and a_2 , i.e., the symbolic value of d is the same when considering either a_1 or a_2 as the duration of the same given tatum at some tempo (respectively $\frac{1}{a_1}$ and $\frac{1}{a_2}$ as shown in Section IV.E). From these definitions, we can now define a *tempo curve* as a path in G . In fact, [9] call such a path a *transcription* rather than a tempo curve, but since an exact tempo curve would be (T_n^*) , those two problems are actually equivalent.

Actually, the consistency property is not that restrictive when considering tempo curves. Let F_1, F_2 be two successive frames, $F_\cap = F_1 \cap F_2$, $d \in F_\cap$, and p a path in G containing a local minima a_1 of ε_{F_1} . According to (11), we can divide the set of all local maxima in two, with those “caused” by F_\cap , M_{F_\cap} , and the others. Let then $m_1, m_2 \in M_{F_\cap}$. Those are local maxima for both ε_{F_1} and ε_{F_2} by (11) since $F_\cap \subset F_2$, and therefore there is at least one local minima within the range $]m_1, m_2[$ for both of these functions. However, thanks to Proposition IV.D.1, we know that both these local minima will quantify the elements of D the same way. Hence, by defining :

- $m_1 = \max\{m \in M_{F_\cap} : m < a_1\}$
- $m_2 = \min\{m \in M_{F_\cap} : m > a_1\}$
- a_2 a local minima of ε_{F_2} in the range $]m_1, m_2[$, which exists since m_1 and m_2 are local maxima (that are not local minima).

We obtain : (a_1, a_2) is *consistent* according to the consistency property.

Corollary IV.D.2 The *consistency property* only implies restrictions relative to the interval of research. In other words, any given strictly partial path p in G can be extended, even if it means considering a bigger interval, for any given performance, and any given frame length for defining G .

However, this restriction on G appears to have some interest. Indeed, let p a path in G locally inconsistent, i.e., such that $a_1, a_2 \in p$ so that $d \in D$ is quantified differently according to a_1 and a_2 , with a_1 and a_2 local minima of successive frames. We therefore have a two partial transcriptions of d being either : m_1 at tempo $\frac{1}{a_1}$ and m_2 at tempo $\frac{1}{a_2}$, m_1, m_2 expressed in tatum unit, with $m_1 \neq m_2$. WHAT IS THE POINT ?

E. Quantification revised

Let us define from now our tatum $\varepsilon = \frac{1}{60}\text{♩}$, which correspond to an sixteenth note wrapped within a triplet within a quintuplet, and has the property that $1 \varepsilon/s = 1 \text{♩}/m$.

With our tatum defined, we can now choose to express all our symbolic units as multiple of this tatum, hence the unit for symbolic values is now ε . We then have : $T = \frac{\Delta b}{\Delta t} = \frac{1\varepsilon}{a}$, where a is the theoretical duration of ε at tempo T . From there, we can define $\sigma_D : a \mapsto \frac{1}{a}\varepsilon_D(a)$, the *normalized error*, or *symbolic error*, since it embodies the error between a transcription of $d \in D$ as m expressed in tatum, hence a quantified and valid transcription, and $d \times \frac{1}{a} = d \times T$, which is the expression of the symbolic duration of d at tempo T according to the [definition of canonical tempo](#).

- LR
- bidi (2 passes: LR + RL) : justification (en annexe) : retour à la définition formelle de Tempo : valide dans les

deux sens, d'où la possibilité de le faire en bidirectionnel + parler rapidement d'une application à Large

- RT : avec valeur initiale de tempo

F. résultats évaluation (comparaison avec 3)

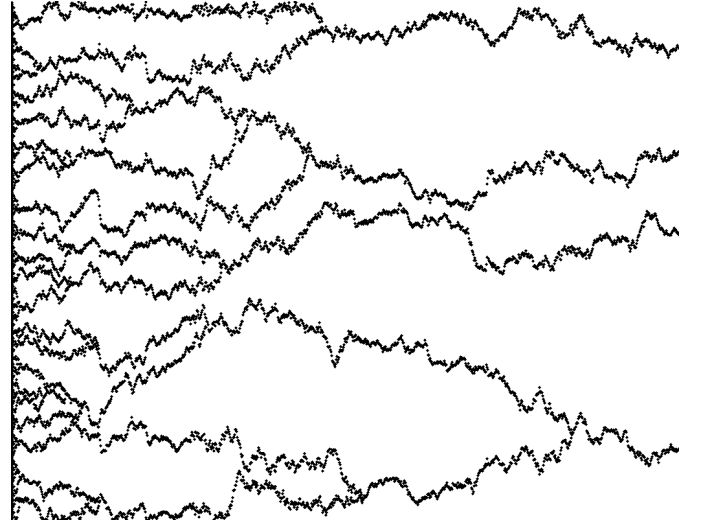


Figure 4: All potentials tempo curves found by a quantified approach for a performance of Piano Sonata No. 11 in A Major (K. 331: III), W.A Mozart. The tempo scale is linear between $\text{♩} = 40$ (bottom) and $\text{♩} = 240$ (top)

V. APPLICATIONS

A. Data extension

Let $(t_n)_{n \in \mathbb{N}}, (T_n^*)_{n \in \mathbb{N}}, (T_n)_{n \in \mathbb{N}}$ be respectively a performance, a canonical tempo for this performance, an estimated tempo curve (supposed to be flattened with respect to the canonical tempo) and $T_c \in R_+^*$ a given tempo value. We define $(s_n := \Delta t_n (T_n - T_n^*))_{n \in \mathbb{N}}$ the symbolic shift of the performance according to (T_n^*) and (T_n) . The *normalized performance* to tempo T_c of (t_n) is then defined as :

$$\hat{t}_0 = t_0, \forall n \in \mathbb{N}, \hat{t}_{n+1} = \hat{t}_n + \underbrace{\Delta t_n \frac{T_n^*}{T_c}}_{\alpha_n} + \underbrace{\frac{s_n}{T_c}}_{\beta_n} \quad (12)$$

Where α_n represents the new duration of Δt_n at tempo T_c , since $\alpha_n = \frac{\Delta b_n}{T_c}$, and β_n embodies the actual time shift at tempo T_c .

Proposition V.A.1 $(\hat{t}_n)_{n \in \mathbb{N}}$ is a performance as defined in Section III.A, and $(\hat{T}_n^*)_{n \in \mathbb{N}} = \left(\frac{1}{1 + \frac{s_n}{\Delta b_n}} T_c \right)_{n \in \mathbb{N}}$

Proof For all $n \in \mathbb{N}^*$, $\Delta \hat{t}_n = \alpha_n + \beta_n = \Delta \frac{\hat{t}_n}{T_c} (T_n^* + T_n - T_n^*) = \Delta \frac{\hat{t}_n}{T_c} T_n > 0$

Furthermore, $\hat{T}_n^* = \frac{\Delta b_n}{\Delta \hat{t}_n} = \frac{\Delta b_n}{\Delta \frac{\hat{t}_n}{T_c} T_n} = \frac{1}{1 + \frac{s_n}{\Delta b_n}} T_c$ \square

Depending on the use of this data, one can choose to adapt the definition of (s_n) , for instance by normalizing its values

or cutting all shifts that represent over a certain portion of their duration. [Proposition V.A.1](#) exposes how such definitions of (s_n) allow for guarantee on tempo change and deviation of the canonical tempo for the *normalized* performance with respect to T_c .

- génération de données “performance” : pour data augmentation ou test robustesse (fuzz testing) aplanissement de tempo démo MIDI?
- transcription MIDI par parsing : pre-processing d’évaluation tempo (approche partie 4)
- analyse “musicologique” quantitative de performances humaines de réf. (à la Mazurka BL) données quantitatives de tempo et time-shifts

VI. CONCLUSION & PERSPECTIVES

- intégration pour couplage avec transcription par parsing (+ plus court chemin multi-critère)
- Section IV.E presented a model which appear to share some similarities with E. W. Large and M. R. Jones [6] as a score-based approach. Therefore, studying the formalism for a quantifier might allow to obtain some theoretical results for the previous model, regarding for instance convergence guarantee and meaningfulness of the result.

REFERENCES

- [1] S. D. Peter, C. E. Cancino-Chacón, E. Karystinaios, and G. Widmer, “Sounding Out Reconstruction Error-Based Evaluation of Generative Models of Expressive Performance,” in *Proceedings of the 10th International Conference on Digital Libraries for Musicology*, 2023, pp. 58–66.
- [2] O. F. B. Katerina Kosta Rafael Ramírez and E. Chew, “Mapping between dynamic markings and performed loudness: a machine learning approach,” *Journal of Mathematics and Music*, vol. 10, no. 2, pp. 149–172, 2016, doi: [10.1080/17459737.2016.1193237](https://doi.org/10.1080/17459737.2016.1193237).
- [3] C. Raphael, “A Probabilistic Expert System for Automatic Musical Accompaniment,” *Journal of Computational and Graphical Statistics*, vol. 10, no. 3, pp. 487–512, Sep. 2001, doi: [10.1198/106186001317115081](https://doi.org/10.1198/106186001317115081).
- [4] “MazurkaBL: Score-aligned Loudness, Beat, and Expressive Markings Data for 2000 Chopin Mazurka Recordings.” Accessed: Jun. 18, 2024. [Online]. Available: <https://zenodo.org/records/1290763>
- [5] P. Hu and G. Widmer, “The Batik-plays-Mozart Corpus: Linking Performance to Score to Musicological Annotations.” Accessed: Jun. 18, 2024. [Online]. Available: <http://arxiv.org/abs/2309.02399>
- [6] E. W. Large and M. R. Jones, “The dynamics of attending: How people track time-varying events,” *Psychological Review*, vol. 106, no. 1, pp. 119–159, 1999, doi: [10.1037/0033-295X.106.1.119](https://doi.org/10.1037/0033-295X.106.1.119).
- [7] H.-H. Schulze, A. Cordes, and D. Vorberg, “Keeping Synchrony While Tempo Changes: Accelerando and Ritardando,” *Music Perception: An Interdisciplinary Journal*, vol. 22, no. 3, pp. 461–477, 2005, doi: [10.1525/mp.2005.22.3.461](https://doi.org/10.1525/mp.2005.22.3.461).
- [8] D. Murphy, “Quantization revisited: a mathematical and computational model,” *Journal of Mathematics and Music*, vol. 5, no. 1, pp. 21–34, Mar. 2011, doi: [10.1080/17459737.2011.573674](https://doi.org/10.1080/17459737.2011.573674).
- [9] G. Romero-García, C. Guichaoua, and E. Chew, “A Model of Rhythm Transcription as Path Selection through Approximate Common Divisor Graphs,” May 2022. Accessed: Jun. 19, 2024. [Online]. Available: <https://hal.science/hal-03714207>
- [10] F. Foscarin, A. Mcleod, P. Rigaux, F. Jacquemard, and M. Sakai, “ASAP: a dataset of aligned scores and performances for piano transcription,” Oct. 2020. Accessed: Jun. 18, 2024. [Online]. Available: <https://cnam.hal.science/hal-02929324>
- [11] S. D. Peter *et al.*, “Automatic Note-Level Score-to-Performance Alignments in the ASAP Dataset,” vol. 6, no. 1, pp. 27–42, Jun. 2023, doi: [10.5334/tismir.149](https://doi.org/10.5334/tismir.149).
- [12] S. Meunier, “Report github,” Jun. 2024, [Online]. Available: <https://github.com/sylvain-meunier/stageL3>
- [13] E. Nakamura, N. Ono, S. Sagayama, and K. Watanabe, “A Stochastic Temporal Model of Polyphonic MIDI Performance with Ornaments,” *Journal of New Music Research*, vol. 44, no. 4, pp. 287–304, Oct. 2015, doi: [10.1080/09298215.2015.1078819](https://doi.org/10.1080/09298215.2015.1078819).
- [14] E. Nakamura, T. Nakamura, Y. Saito, N. Ono, and S. Sagayama, “Outer-Product Hidden Markov Model and Polyphonic MIDI Score Following,” *Journal of New Music Research*, vol. 43, no. 2, pp. 183–201, Apr. 2014, doi: [10.1080/09298215.2014.884145](https://doi.org/10.1080/09298215.2014.884145).
- [15] K. Shibata, E. Nakamura, and K. Yoshii, “Non-local musical statistics as guides for audio-to-score piano transcription,” *Information Sciences*, vol. 566, pp. 262–280, Aug. 2021, doi: [10.1016/j.ins.2021.03.014](https://doi.org/10.1016/j.ins.2021.03.014).
- [16] J. Hentschel, M. Neuwirth, and M. Rohrmeier, “The Annotated Mozart Sonatas: Score, Harmony, and Ca-

dence,” vol. 4, no. 1, pp. 67–80, May 2021, doi: [10.5334/tismir.63](https://doi.org/10.5334/tismir.63).

- [17] M. Müller, “Memory-restricted Multiscale Dynamic Time Warping,” [Online]. Available: https://www.academia.edu/25724042/MEMORY_RESTRICTED_MULTISCALE_DYNAMIC_TIME_WARPING
- [18] E. Nakamura, K. Yoshii, and H. Katayose, “Performance Error Detection and Post-Processing for Fast and Accurate Symbolic Music Alignment,” 2017. Accessed: Jun. 18, 2024. [Online]. Available: <https://www.semanticscholar.org/paper/Performance-Error-Detection-and-Post-Processing-for-Nakamura-Yoshii/37e9f5e23cada918c2b8982d71a18972140d9d5a>
- [19] A. Cont, “ANTESCOFO: Anticipatory Synchronization and Control of Interactive Parameters in Computer Music,” in *International Computer Music Conference (ICMC)*, 2008, pp. 33–40.
- [20] E. W. Large *et al.*, “Dynamic models for musical rhythm perception and coordination,” *Frontiers in Computational Neuroscience*, vol. 17, May 2023, doi: [10.3389/fncom.2023.1151895](https://doi.org/10.3389/fncom.2023.1151895).
- [21] J. D. Loehr, E. W. Large, and C. Palmer, “Temporal coordination and adaptation to rate change in music performance,” *Journal of Experimental Psychology. Human Perception and Performance*, vol. 37, no. 4, pp. 1292–1309, Aug. 2011, doi: [10.1037/a0023102](https://doi.org/10.1037/a0023102).

APPENDIX A : SOME FORMAL CONSIDERATIONS

A. Equivalence of tempo formal definitions

Let $n \in \mathbb{N}$. $\int_{t_0}^{t_n} T(t) dt = \sum_{i=0}^{n-1} \int_{t_i}^{t_{i+1}} T(t) dt$
 Furthermore, $\int_{t_0}^{t_{n+1}} T(t) dt = \int_{t_0}^{t_{n+1}} T(t) dt + \int_{t_n}^{t_0} T(t) dt = \int_{t_0}^{t_{n+1}} T(t) dt - \int_{t_0}^{t_n} T(t) dt$.

Let T be a formal tempo according to the first definition. For all $n \in \mathbb{N}$, we then have :

$$\begin{aligned} \int_{t_n}^{t_{n+1}} T(t) dt &= \int_{t_0}^{t_{n+1}} T(t) dt - \int_{t_0}^{t_n} T(t) dt \\ &= b_{n+1} - b_0 - (b_n - b_0) \\ &= b_{n+1} - b_n \end{aligned}$$

Let T be a formal tempo according to the second definition.

For all $n \in \mathbb{N}$:

$$\begin{aligned} \int_{t_0}^{t_n} T(t) dt &= \sum_{i=0}^{n-1} \int_{t_i}^{t_{i+1}} T(t) dt \\ &= \sum_{i=0}^{n-1} b_{i+1} - b_i \\ &= b_n - b_0 \end{aligned}$$

We thus obtain the two implications, hence the equivalence.

B. The tempo octave problem

When estimating tempo, or transcribing a performance, there always exist several equivalent possibilities. For instance, given a “correct” transcription (b_n) of a performance (t_n) , one can choose to define its own transcription as $t = \left(\frac{b_n}{2}\right)_{n \in \mathbb{N}}$.

Then, the canonical tempo according to t , called (T_1^*) , and the one according to (b_n) , called (T_2^*) verify :

$\forall n \in \mathbb{N}, T_{1,n}^* = \frac{\frac{b_{n+1}}{2} - \frac{b_n}{2}}{t_{n+1} - t_n} = \frac{1}{2} T_{2,n}^*$. Actually, the t transcription corrections to (b_n) where all durations are indicated doubled, but played twice faster, hence giving the exact same theoretical performance. Unfortunately, there is no absolute way to decide which of those two transcription is better than the other. This problem is here known as the tempo octave problem, and should be kept in mind when transcribing or estimating tempo. We present in Section IV.B a model resistant to these tempo octaves, as well as other kind of octaves not discussed here (for instance multiplying the tempo by 3 by using [triplet](#)).

C. Tempo conservation when reversing time

First, we want to insist on the fact that none of the sequence (b_n) and (t_n) are infinite, but in order to simplify the notation, we chose to indicate them as usual infinite sequences, or rather only consider them on a finite number of indexes, called $|(b_n)|$ and $|(t_n)|$ respectively, with $|(b_n)| = |(t_n)|$. Let us then define the reversed sequence of $(u_n)_{n \in \mathbb{N}}$ as $r((u_n)_{n \in \mathbb{N}}) := (\bar{u}_n = u_{|(u_n)| - n})_{n \in \llbracket 0, |(u_n)| \rrbracket}$. Both $\bar{b} = r((b_n)_{n \in \mathbb{N}})$

and $\bar{t} = r((t_n)_{n \in \mathbb{N}})$ are correct representations of a sheet music and performance respectively, as defined in Section I-II.A.

Let $t^* = t_{|(t_n)|}$ and $q = |(t_n)|$, T a formal tempo with respect to (t_n) and (b_n) , i.e., $\forall n \in \mathbb{N}$, $\int_{t_n}^{t_{n+1}} T(t) dt = b_{n+1} - b_n$, $n \in \llbracket 0, q-1 \rrbracket$, and $T_r : t \mapsto T(t^* - t)$.

$$\begin{aligned} \int_{t_n}^{t_{n+1}} T_r(t) dt &= \int_{t^*-t_{q-n-1}}^{t^*-t_{q-n}} T(t^* - t) dt = \int_{t_{q-n}}^{t_{q-n-1}} -T(x) dx \\ &= \int_{t_{q-n-1}}^{t_{q-n}} T(t) dt \\ &= b_{q-n} - b_{q-n-1} \\ &= (b_{q-n} - b_q) - (b_{q-n-1} - b_q) \\ &= -\bar{b}_n + \bar{b}_{n+1} \\ &= \bar{b}_{n+1} - \bar{b}_n \end{aligned}$$

Hence T_r is a formal tempo with respect to (\bar{t}_n) and (\bar{b}_n) .

D. Musical explication of the choice of a tempo distance

In terms of tempo, halving and doubling are considered as far as each other from the initial value. Therefore a usual absolute distance does not fit this notion, and we will rather use a logarithmic distance when comparing tempi.

APPENDIX B : ESTIMATOR MODEL

A. Formal explanations and proofs

First d is indeed a mathematical distance : let $a, b \in (R_+^*)^2$, $d(a, b) = d(b, a)$ and $d(a, b) = 0 \Leftrightarrow |\log(\frac{a}{b})| = 0 \Leftrightarrow a = b$. Finally, let $c \in R_+^*$, $d(a, c) = k_* |\log(\frac{a}{c})| = k_* |\log(\frac{a}{b} \times \frac{b}{c})| = k_* |\log(\frac{a}{b}) + \log(\frac{b}{c})| \leq d(a, b) + d(b, c)$.

Then, (8) presents an argmin, that makes sense when E is a increasing right-continuous function-like object, even though its actual expression may change after each computed value of T_{n+1} . In fact, E can only output a countable set of values, hence E is piecewise constant under those hypothesis.

Finally, one can notice that the value of $k_* \in R_+^*$ does not affect the result of the process.

B. About the range $[\frac{\sqrt{2}}{2}T_n, \sqrt{2}T_n]$

In order to resist to the tempo octave problem discussed in Appendix A, we choose here to consider a unique candidate within a range $[x, 2x] \subset [\frac{1}{2}, 2]$, for a given $x \in R_+^*$. Then, we want this range to be centered around 1, since its values corresponds to tempo variation, and our system should not favor increasing nor decreasing the tempo *a priori*. For this musical reason, we then take x as solution of : $\|x - 1\| = \|2x - 1\|$ that implies $1 - x = 2x - 1$, i.e., $x = \frac{2}{3}$ with the absolute value distance.

With a logarithmic distance, the same reasoning would give : $\log(\frac{1}{x}) = \log(2x) \Leftrightarrow -\log(x) = \log(2) + \log(x) \Leftrightarrow \log(x^2) = -\log(2) \Leftrightarrow x^2 = \frac{1}{2} \Leftrightarrow x = \frac{\sqrt{2}}{2}$ since $x > 0$.

Then, when considering the tempo distance between T_{n+1} and T_n , we find :

$$d(T_{n+1}, T_n) = k_* |\log(y)| = d(1, y) \quad (13)$$

where $y = \operatorname{argmin}_{x' \in [x, 2x]} d\left(x', \frac{\Delta t_n}{\Delta t_{n+1}} E\left(x' \frac{\Delta t_{n+1}}{\Delta t_n}\right)\right)$.

Therefore, since we want the extreme possible values of our range to imply an equal distance between T_n and T_{n+1} , we choose the logarithmic distance, and hence $x = \frac{\sqrt{2}}{2}$, so that $d(1, x) = d(1, 2x)$.

C. About the estimator E

One can notice that $E = \text{id}$ implies, by the hypothesis that E acts as an oracle, that the theoretical and actual values are the same, or that the performance is a perfect interpretation of the piece. Since real players do not make such performance, we can expect a relevant estimator to act rather differently than the identity function.

Moreover, E is not a function : its expression only has to be fixed when computing the numerical resolution for the argmin. Hence, an given output can depends on several previous outputs. In an extreme case, E can even be a transcripting system. However, in our problem of tempo estimation, we do not have as much constraints as in transcription.

Indeed, the following figures displays two transcription A and B, the latter being incorrect with regards to usual transcription convention, and their corresponding tempo curves.

One can notice that these are quite similar, and in fact, a human being could not tell them apart, as shown by Figure 8.



Figure 5: Transcription A

Figure 6: Transcription B

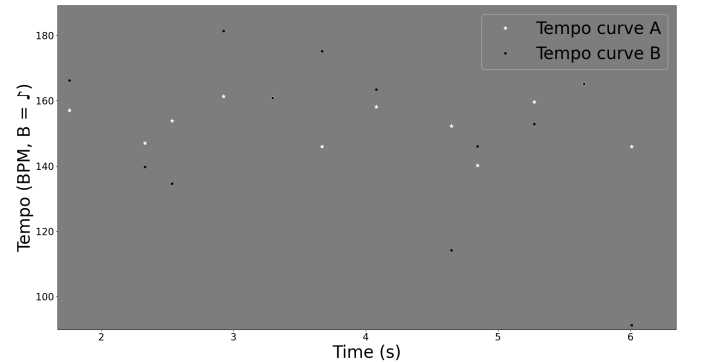


Figure 7: Tempo curves A & B

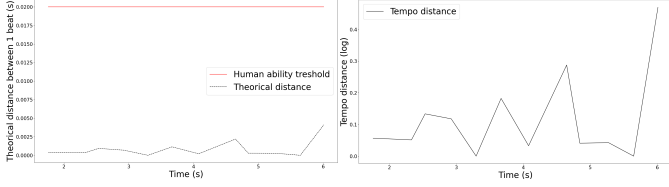


Figure 8: Tempo distance (s) Figure 9: Tempo distance (log)

Tempo distance between the two previous curves. Being able to differentiate them would imply to tell apart two rhythmic events within 4 ms, which is suppose impossible for a humain being according to the value of ε defined (and displayed as the top line in Figure 8) in Section III.A.

APPENDIX C : QUANTIFIED MODEL

In this section we use the notation introduced by G. Romero-García, C. Guichaoua, and E. Chew [9], that essentially replace D by T .

Let us first define : $g : x \mapsto \min(x - \lfloor x \rfloor, 1 + \lfloor x \rfloor - x)$

One can make sure that $g : x \mapsto \begin{cases} x - \lfloor x \rfloor & \text{si } x - \lfloor x \rfloor \leq \frac{1}{2} \\ 1 - (x - \lfloor x \rfloor) & \text{sinon} \end{cases}$ and that g is 1-periodic, continuous on \mathbb{R} .

Then, by definition :

$$\begin{aligned} \varepsilon_T(a) &= \max_{t \in T} \left(\min_{m \in \mathbb{Z}} |t - ma| \right) \\ &= \max_{t \in T} \min(t - \lfloor \frac{t}{a} \rfloor a, (\lfloor \frac{t}{a} \rfloor + 1)a - t) \\ &= a \max_{t \in T} \min \left(\frac{t}{a} - \lfloor \frac{t}{a} \rfloor, \lfloor \frac{t}{a} \rfloor + 1 - \frac{t}{a} \right) \\ &\quad \underbrace{\hspace{10em}}_{g(\frac{t}{a})} \\ &= a \max_{t \in T} g\left(\frac{t}{a}\right) \end{aligned}$$

hence we have proved ε_T to be continous on R_+^* .

Furthermore, for $n \in \mathbb{N}^*, T \subset (\mathbb{R}_+^*)^n, a \in R_+^*$,

$$\varepsilon_T(a) = a \max_{t \in T} g\left(\frac{t}{a}\right) = a \times \max_{t \in T} g\left(\frac{t}{a} 1^{-1}\right) = a \varepsilon_{T/a}(1).$$

Hence the intuitive following result : the smaller the tatum, the smaller the bound of the error.

A. Characterization of local maxima that happen not be a local minima

1) *First implication:*

Let a be a local maxima (and not a local minima) of $\varepsilon_T, a > 0$. By definition, there is a $a > \varepsilon > 0$ so that :

$$\forall \delta \in]-\varepsilon, \varepsilon[, \varepsilon_T(a) \geq \varepsilon_T(a + \delta).$$

Let $t \in \operatorname{argmax}_{t' \in T} g\left(\frac{t'}{a}\right)$, hence $\varepsilon_T(a) = ag\left(\frac{t}{a}\right)$.

For all $\delta \in]-\varepsilon, \varepsilon[, \varepsilon_T(a) = ag\left(\frac{t}{a}\right) \geq \varepsilon_T(a + \delta)$,

$$\text{and } \varepsilon_T(a + \delta) = (a + \delta) \max_{t' \in T} g\left(\frac{t'}{a + \delta}\right) \geq (a + \delta) g\left(\frac{t}{a + \delta}\right).$$

For $\delta \geq 0, a + \delta \geq a$, so $ag\left(\frac{t}{a}\right) \geq (a + \delta)g\left(\frac{t}{a + \delta}\right) \geq ag\left(\frac{t}{a + \delta}\right)$. Hence, $g\left(\frac{t}{a}\right) \geq g\left(\frac{t}{a + \delta}\right)$ since $a > 0$, for all $\delta \in [0, \varepsilon]$.

Therefore, g increases monotonically in the range $]\frac{t}{a + \delta}, \frac{t}{a}[$, since g has a unique local maxima (modulo 1), considering the previous range as a neighbourhood of $\frac{t}{a}$.

Hence, $g = x \mapsto x - \lfloor x \rfloor$ within the considered range, and $g\left(\frac{t}{a}\right) = \frac{t}{a} - \lfloor \frac{t}{a} \rfloor, \varepsilon_T(a) = t - a \lfloor \frac{t}{a} \rfloor$.

The function ε_T is the maximum of a finite set of continuous functions, with a countable set of A of intersection, i.e., $A = \{x \in \mathbb{R}_+^* : \exists (t_1, t_2) \in T^2 : t_1 \neq t_2 \wedge xg\left(\frac{t_1}{x}\right) = xg\left(\frac{t_2}{x}\right)\}$. Indeed,

$$x \in A \Leftrightarrow \exists (t_1, t_2) \in T^2 : t_1 \neq t_2 \wedge g\left(\frac{t_1}{x}\right) = g\left(\frac{t_2}{x}\right)$$

$$\Leftrightarrow \exists (t_1, t_2) \in T^2 : t_1 \neq t_2 \wedge \frac{t_1}{x} = \pm \frac{t_2}{x} \pmod{1}$$

$$\Leftrightarrow \exists (t_1, t_2) \in T^2 : t_1 \neq t_2 \wedge x = \frac{t_1 \mp t_2}{n}, n \in \mathbb{Z}^*$$

Hence $A \subset \left\{ \frac{t_1 \mp t_2}{n}, (t_1, t_2) \in T^2, n \in \mathbb{Z}^* \right\}$, because $T \subset (R_+^*)^{|T|}$. Therefore, there is a countable set of closed convex intervals, whose union is R_+^* so that on each of these intervals, ε_T is equal to $f_t : a \mapsto ag\left(\frac{t}{a}\right)$ for a $t \in T$. Let then t be so that for all $x \in]a - \delta', a]$, $\varepsilon_T(x) = f_t(x)$, where $]a - \delta', a]$ is included in one the previous intervals. Since f_t and ε_T are both continuous on $]a - \delta', a]$, $f_t(a) = \varepsilon_T(a)$ and therefore, $t \in \operatorname{argmax}_{t' \in T} g\left(\frac{t'}{a}\right)$. The previous paragraph showed that g is increasing on a left neighbourhood of $\frac{t}{a}$. Therefore, on a right neighbourhood of $\frac{t}{a}$, g is either increasing or decreasing by its definition.

- if g is increasing on this neighbourhood, called $N\left(\frac{t}{a}\right)^+$ in the following, the previous expression of g remains valid, i.e., $\forall x \in N\left(\frac{t}{a}\right)^+, g(x) = x - \lfloor x \rfloor$. Moreover, $x \mapsto \lfloor x \rfloor$ is right-continuous, hence by restricting $N\left(\frac{t}{a}\right)^+$, we can assure for all $x \in N\left(\frac{t}{a}\right)^+, \lfloor x \rfloor = \lfloor \frac{t}{a} \rfloor$. Let then $y = a - \frac{t}{x}$ so that $x = \frac{t}{a - y}$, we then have $\varepsilon_T(a - y) = t - (a - y) \lfloor \frac{t}{a} \rfloor \leq \varepsilon_T(a) = t - a \lfloor \frac{t}{a} \rfloor$ because a is a local maxima of ε_T and $a - y$ is within a (left) neighbourhood of a , even if it means restricting δ' or $N\left(\frac{t}{a}\right)^+$. Hence, $t - \lfloor \frac{t}{a} \rfloor a \geq t - (a - y) \lfloor \frac{t}{a} \rfloor$ i.e., $a \lfloor \frac{t}{a} \rfloor \leq (a - y) \lfloor \frac{t}{a} \rfloor \Leftrightarrow 0 \leq -y \lfloor \frac{t}{a} \rfloor$ i.e., $\lfloor \frac{t}{a} \rfloor \leq 0$ i.e., $\lfloor \frac{t}{a} \rfloor = 0$, since y, t and a are all positive values. Then, $a > t$ and therefore $\varepsilon_T(a - y) = t = \varepsilon_T(a)$ for $\lfloor \frac{t}{a} \rfloor = 0$. This interval where ε_T is constant is then either going on infinitely on the right of a , or else ε_T will reach a value greater than $\varepsilon_T(a) = t$, since ε_T can then be rewritten as $x \mapsto \max \left(\max_{t' \in T \setminus \{t\}} xg\left(\frac{t'}{x}\right), \varepsilon_T(a) \right)$ on $[a, +\infty[$. Hence a is a local minima on the right, and since ε_T is constant on a left neighbourhood of a , a is also a local minima on the left. Finally, a is a local minima, which is absurd by definition.
- else, g is decreasing on $N\left(\frac{t}{a}\right)^+$, $\frac{t}{a}$ is by definition a local maxima of g . However, g only has a unique local maxima modulo 1, that is $\frac{1}{2}$. Hence, $\frac{t}{a} = \frac{1}{2} \pmod{1}$, i.e., $\frac{t}{a} = \frac{1}{2} + k, k \in \mathbb{Z}$, or $a = \frac{t}{\frac{1}{2} + k}, k \in \mathbb{N}$, since $a > 0$.

2) Second implication:

Let $(t, k) \in T \times \mathbb{N}$, $a = \frac{t}{k + \frac{1}{2}}$.

By definition: $g(\frac{t}{a}) = g(\frac{1}{\frac{1}{2} + k}) = g(\frac{1}{2}) = \frac{1}{2} = \max_{\mathbb{R}} g$.

Therefore, $\varepsilon_T(a) = a \max_{t' \in T} g(\frac{t'}{a}) = ag(\frac{t}{a}) = \frac{a}{2}$.

For all $x \in]0, a[$, $\varepsilon_T(x) = x \max_{t' \in T} g(\frac{t'}{x}) \leq \frac{x}{2} < \frac{a}{2} = \varepsilon_T(a)$, i.e., $\varepsilon_T(a) > \varepsilon_T(x)$.

Let $T^* = \{t' \in T : g(\frac{t'}{a}) = \frac{1}{2}\}$. Since $t \in T^*$, $|T^*| > 1$.

Let $t^* \in T^*$. For all $t' \in T \setminus T^*$, $g(\frac{t'}{a}) < g(\frac{t^*}{a})$.

Since $h_{t'} : x \mapsto g(\frac{t'}{x}) - g(\frac{t^*}{x})$ is continuous in a neighbourhood of $a > 0$, we have the existence of $\varepsilon_{t'} > 0$ so that $h_{t'}$ is strictly positive within $[a, a + \varepsilon_{t'}]$.

Let $\varepsilon_{t^*} = \min_{t' \in T} \varepsilon_{t'}$ and finally $\varepsilon_1 = \min_{t^* \in T^*} \varepsilon_{t^*}$.

Let $(t_1, t_2) \in (T^*)^2$.

In the following, $N(a)^+$ is a right neighbourhood of a such that $a \notin N(a)^+$.

Let $\text{tmp} : x \mapsto g(\frac{t_1}{x}) - g(\frac{t_2}{x})$ be a continuous function on $N(a)^+$ and A be the set of all $x^* \in N(a)^+$ so that $\text{tmp}(x^*) = 0 \Leftrightarrow g(\frac{t_1}{x^*}) = g(\frac{t_2}{x^*})$.

We have for all $x^* \in A$, $g(\frac{t_1}{x^*}) = g(\frac{t_2}{x^*})$ by definition. Considering the expression of g , we then find: $\frac{t_1}{x^*} = \pm \frac{t_2}{x^*} \bmod 1$. Moreover, since g only reach $g(\frac{t_1}{a}) = \frac{1}{2}$ once per period, we have $\frac{t_1}{a} = \frac{t_2}{a} \bmod 1$, i.e., $|\frac{t_1}{a} - \frac{t_2}{a}| = k_a \in \mathbb{N}$.

Then, $\frac{t_1}{x^*} = \pm \frac{t_2}{x^*} \bmod 1$ i.e., $|\frac{t_1}{x^*} \mp \frac{t_2}{x^*}| = k_* \in \mathbb{N}$, and therefore $|t_1 \mp t_2| = ak_* = x^* k_*$, and $x^* > a$ implies $k_a > k_* \geq 0$. However, $x^* = \frac{|t_1 \mp t_2|}{k_*}$, hence A is finite if $A \neq \emptyset$, and \emptyset is a finite set. Finally, A is a finite set, i.e., $|A| \in \mathbb{N}$.

Let then $x_{t_1, t_2} = \begin{cases} \min A \text{ if } A \neq \emptyset \\ x \in N(a)^+ \setminus \{a\} \text{ otherwise} \end{cases}$ **WLOG**, $g(\frac{t_1}{x}) \geq g(\frac{t_2}{x}) \forall x \in [a, x_{t_1, t_2}]$

Let $a_2 = \min_{(t_1, t_2) \in T^{*2}} x_{t_1, t_2}$ and $a_1 \in]a, a_2[$,

let $t^* = \arg\max_{t' \in T^*} g(\frac{t'}{a_1})$ We finally have $\forall x \in]a, a_2[$, $g(\frac{t^*}{x}) \geq g(\frac{t'}{x})$, $\forall t' \in T^*$.

Let then $\tilde{a} = \min(a + \varepsilon_1, a_2)$ so that for all $x \in]a, \tilde{a}[$, $t' \in T$, $g(\frac{t^*}{x}) \geq g(\frac{t'}{x})$, hence $\varepsilon_T(x) = xg(\frac{t^*}{x})$.

Let $f : x \mapsto g(\frac{t^*}{x})$, $f(a) = g(\frac{t^*}{a}) = \frac{1}{2}$ because $t^* \in T^*$ hence f is increasing on a right neighbourhood of a , $N(a)^+$, since $\frac{1}{2}$ is a global maxima of g , therefore g is increasing on $N(\frac{t^*}{a})$ a left neighbourhood of $\frac{t^*}{a}$, i.e., f is increasing on $N(a)^+$. Therefore, we know that $g(\frac{t^*}{x}) = \frac{t^*}{x} - \lfloor \frac{t^*}{x} \rfloor$, since the only other possible expression for g would imply a decreasing function on $N(\frac{t^*}{a})^-$.

Hence, $\varepsilon_T(x) = xg(\frac{t^*}{x}) = x(\frac{t^*}{x} - \lfloor \frac{t^*}{x} \rfloor) = t^* - x \lfloor \frac{t^*}{x} \rfloor$ and $\varepsilon_T(a) = t^* - a \lfloor \frac{t^*}{a} \rfloor$ since ε_T is continuous on \mathbb{R}_+^* .

By definition: $\lfloor \frac{t^*}{a} \rfloor \leq \frac{t^*}{a} < \lfloor \frac{t^*}{a} \rfloor + 1$.

However, $f(a) = \frac{1}{2} = \frac{t^*}{a} - \lfloor \frac{t^*}{a} \rfloor$, therefore $\lfloor \frac{t^*}{a} \rfloor < \frac{t^*}{a}$.

Then, there is $\alpha \in \mathbb{R}_+^*$ so that $\lfloor \frac{t^*}{a} \rfloor < \alpha < \frac{t^*}{a}$, let $y = \frac{t^*}{\alpha}$, i.e.,

$\alpha = \frac{t^*}{y}$, with $y > a$.

Let $a' = \min(y, \tilde{a})$ and $k = \lfloor \frac{t^*}{a'} \rfloor$. For all $x \in]a, a'[$,

- $x \leq y = \frac{t^*}{\alpha}$ hence $\lfloor \frac{t^*}{a'} \rfloor \leq \alpha \leq \frac{t^*}{x}$
- $x \geq a$ hence $\frac{t^*}{x} \leq \frac{t^*}{a} < \lfloor \frac{t^*}{a} \rfloor + 1$

In the end, $\lfloor \frac{t^*}{x} \rfloor = \lfloor \frac{t^*}{a'} \rfloor = k$ by definition.

Then, $\varepsilon_T(x) = t^* - xk$ and $\varepsilon_T(a) = t^* - xa$, with $a < x$.

Finally, $\varepsilon_T(a) > \varepsilon_T(x)$.

To conclude, for all $x \in]0, a'[$,

- if $x \leq a$, $\varepsilon_T(a) \geq \varepsilon_T(x)$
- if $x \geq a$, $x \in [a, a'[$ and $\varepsilon_T(a) \geq \varepsilon_T(x)$

Hence a is a local maxima of ε_T and not a local minima, and then the set of all such local maxima of ε_T is $M_T = \left\{ \frac{t}{k + \frac{1}{2}}, t \in T, k \in \mathbb{N} \right\}$ and therefore with the notations introduced in Section IV.D, we proved (11) ■

B. Necessary condition to be a local minima

Let a be a local minima of ε_T , i.e.,

ε_T is constant on a neighbourhood of $+\infty$, donc on va d'abord trouver un maximum local qui n'est pas un minimum local, puis le reste de la preuve fonctionne

Par continuité de ε_T , on est assuré de l'existence d'exactly un unique minimum local entre deux maximums locaux, qui est alors global sur cet intervalle.

Par la condition nécessaire précédente, il suffit donc, pour déterminer ce minimum local, de déterminer le plus petit élément parmi les points obtenus, contenus dans l'intervalle. On en déduit ainsi un algorithme en $\mathcal{O}(|T|^2 \frac{t^*}{\tau} \log(|T| \frac{t^*}{\tau}))$ permettant de déterminer tous les minimums locaux accordés par le seuil τ fixé, sur l'intervalle $]2\tau, t_* + \tau[$

On en déduit la correction de Algorithm 1.

APPENDIX D : MUSICAL GLOSSARY

A. Acronyms

MIR – Music Information Retrieval: Interdisciplinary science aiming at retrieving information from music, in several ways. Amongst the various problems tackled by the community, one can notice [transcription](#), automatic or semi-automatic musical analysis, performance generation or classification... 1

MTU – Musical Time Unit: Time unit for a symbolic, or musical notation, e.g., beat, quarter note (♩), eighth note (♪). 1

RTU – Real Time Unit: Time unit to represent real events.

Here, we usually used seconds as RTU. 1

WLOG – Without loss of generality: The term is used to indicate the assumption that what follows is chosen arbitrarily, narrowing the premise to a particular case, but does not affect the validity of the proof in general. The other

cases are sufficiently similar to the one presented that proving them follows by essentially the same logic. 12

B. Definitions

articulation: Describes how a specific note is played by the performer. For instance, *staccato* means the note shall not be maintained, and instead last only a few musical units, depending on the context. On the other hand, a fermata (*point d'orgue* in French) indicates that the note should stay longer than indicated, to the performer's discretion. 1

beat: Symbolic time unit of a score, its value is defined by a time signature. Although its value can change within a score, or through various transcription of a same piece, this notion is usually the most convenient way to describe a rhythmic sequence of events. 1, 13

chord: A chord is by definition the simultaneous production of at least three musical events with different pitches 2

measure: A measure is a symbolic time unit corresponding to a fixed amount (integer) of beats. This value is indicated by the *time signature*. 2

online: In computer science, an online algorithm is one that can process its input piece-by-piece in a serial fashion, i.e., in the order that the input is fed to the algorithm, without having the entire input available from the start. 5

rest: A symbolic notation for silence, following the same rules as actual note notations. 2

tatum: Minimal resolution of a musical unit, expressed in beats. Although several values are possible, a tatum is usually indicates the following value for a given score (b_n) : $\sup\{r \mid \forall n \in \mathbb{N}, \exists k \in \mathbb{N} : b_n = kr, r \in \mathbb{R}_+^*\}$. For practical reasons, a tatum may be defined as smaller value than the one previously given, especially if this value is easier to express within the current time signature, or makes more sense musically.. 4, 6

tempo: Formally defined in Section III by : $T_n^* = \frac{b_{n+1}-b_n}{t_{n+1}-t_n}$, tempo is a measure of the immediate speed of a performance, usually written on the score. It can be seen as a ratio between the symbolic speed indicated by the score, and the actual speed of a performance. Tempo is often expressed in *beat* per minute, or bpm 1

time signature: The time signature is a convention in Western music notation that specifies how many note values of a particular type are contained in each measure. It is composed of two integers : the amount of beat contained within a measure, and the value of these beats, indicated as division of a whole note, i.e., four quarter notes. 13

transcription: Process of converting an audio recording into symbolic notation, such as sheet music or MIDI file. This process involves several audio analysis tasks, which may include multi-pitch detection, duration or tempo estimation, instrument identification... 2, 12

triplet: A triplet is a musical symbol indicating to play a third of the indicated duration. 9

velocity: The velocity describes how loud a sound shall be played, or is actually played. 1