

# An Analysis Method for Metric-Level Switching in Beat Tracking

Ching-Yu Chiu, Meinard Müller, Matthew E. P. Davies, Alvin Wen-Yu Su, and Yi-Hsuan Yang

**Abstract**—For expressive music, the tempo may change over time, posing challenges to tracking the beats by an automatic model. The model may first tap to the correct tempo, but then may fail to adapt to a tempo change, or switch between several incorrect but perceptually plausible ones (e.g., half- or double-tempo). Existing evaluation metrics for beat tracking do not reflect such behaviors, as they typically assume a fixed relationship between the reference beats and estimated beats. In this paper, we propose a new performance analysis method, called annotation coverage ratio (ACR), that accounts for a variety of possible metric-level switching behaviors of beat trackers. The idea is to derive sequences of modified reference beats of all metrical levels for every two consecutive reference beats, and compare every sequence of modified reference beats to the subsequences of estimated beats. We show via experiments on three datasets of different genres the usefulness of ACR when being utilized alongside existing metrics, and discuss the new insights that can be gained.

**Index Terms**—Beat tracking, evaluation metrics

## I. INTRODUCTION

**B**EAT tracking aims to automatically determine a sequence of time positions that a listener would tap to when listening to a piece of music [1]–[3]. Most approaches, which have been found to work well for tracking the beats of music with steady tempo [4]–[10], follow a two-block architecture: use a deep learning-based block to generate continuous-valued beat activation functions  $\Delta(n)$  indicating the likelihood of observing a beat at a time instance  $n$ , and a post processing tracker (PPT) that determines the final beat positions from  $\Delta(n)$  using, for example, dynamic programming (DP) [3], [11] or a hidden Markov model (HMM) [12], [13]. See Fig. 1(a).

To quantify the performance of a beat tracker and to study the avenues for improvement, an objective evaluation metric that compares the estimated beat positions  $\hat{B} = \{\hat{b}_1, \hat{b}_2, \dots\}$

Meinard Müller is supported by a grant from the International Audio Laboratories Erlangen, which are a joint institution of the Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU) and Fraunhofer Institute for Integrated Circuits IIS. Matthew E. P. Davies is funded by national funds through the FCT - Foundation for Science and Technology, I.P., within the scope of the project CISUC - UID/CEC/00326/2020 and by European Social Fund, through the Regional Operational Program Centro 2020.

Ching-Yu Chiu is with the Graduate Program of Multimedia Systems and Intelligent Computing, National Cheng Kung University and Academia Sinica, Taiwan (e-mail: sunnycyc@citi.sinica.edu.tw).

Meinard Müller is with the International Audio Laboratories Erlangen, Germany (e-mail: meinard.mueller@audiolabs-erlangen.de).

Matthew E. P. Davies is with the Department of Informatics Engineering, Centre for Informatics and Systems of the University of Coimbra, University of Coimbra, Portugal (e-mail: mepdavies@dei.uc.pt).

Alvin Wen-Yu Su is with the Department of CSIE, National Cheng Kung University, Taiwan (e-mail: alvinsu@mail.ncku.edu.tw).

Yi-Hsuan Yang is with Taiwan AI Labs. (e-mail: yhyang@ailabs.tw).

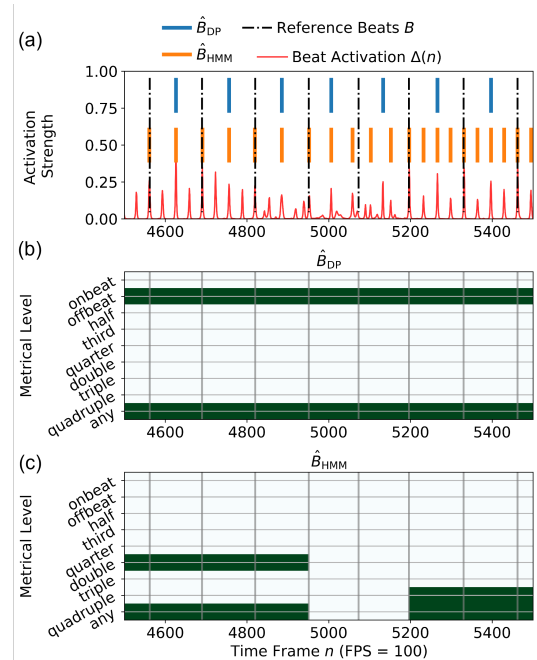


Fig. 1. Visualization of the ACR ( $L = 2$ ) evaluation showing the Boolean values whether a reference beat is “L-correct detected” by the estimated beats w.r.t. different metrical levels. (a) Input beat activation function  $\Delta(n)$ , reference beats  $B$  and the beat estimation of two trackers (DP HMM). Visualization of (b) the DP-based and (c) the HMM-based result.

with the reference (ground truth) ones  $B = \{b_1, b_2, \dots\}$  is needed. As the focus has been mostly on tracking music with steady tempo, most existing metrics assume that the relationship between  $B$  and  $\hat{B}$  with regard to their metrical level is fixed [14]–[16]. For example,  $B$  and  $\hat{B}$  might be tapped at the same tempo, or  $\hat{B}$  might be tapped at a different yet plausible tempo (e.g., half or double tempo) throughout a piece. Accordingly, if a beat tracker switches its metrical level in the middle of a piece, such metrics may underestimate the beat tracker’s performance. For example, the HMM tracker in Fig. 1(a) taps at double tempo first but switches to quadruple tempo later; existing metrics would only consider its result to be mostly incorrect, while perceptually a listener may regard it as performing well most of the time [17]–[22]. Such a metric-level switching (MLS) issue has been identified by researchers before [17], [23], [24], but there is relatively little work on explicitly investigating MLS.<sup>1</sup>

<sup>1</sup>Strictly speaking, “metric-level” considers only different tempi. However, as tapping on offbeats (i.e., phases) can be a common mistake (or behavior) of beat trackers, we follow [16], [25] and consider a broader definition (with slight abuse of language) that includes offbeats as a type of MLS.

TABLE I  
 CONDITIONS CONSIDERED BY DIFFERENT METRICS; ‘(✓)’ DENOTES  
 ‘PARTIALLY’ (CF. SECTION II-C)

Conditions	Evaluation metrics			
	F1	CMLt/AMLt	L-correct	ACR
<i>onbeat</i>	✓	✓	✓	✓
<i>offbeat</i> (half)		✓	✓	✓
<i>offbeat</i> (one-third, two-third)			✓	✓
<i>subharmonic</i> (half, one-third)		✓		✓
<i>subharmonic</i> (quarter)				✓
<i>harmonic</i> (double, triple)		✓		✓
<i>harmonic</i> (quadruple)				✓
<i>metric-level switching</i> (MLS)			(✓)	✓

According to our previous attempts to build beat trackers for expressive music that features rich tempo variations [26], such an MLS issue is prevalent, making it hard to adequately compare the performance of different approaches. This motivates us to propose a new metric.

In this paper, we propose a novel method called “annotation coverage ratio” (ACR) that addresses MLS and considers various common metrical levels while evaluating the correctness of each estimated beat. As listed in Table I, ACR considers a wider variety of metric-level relationships between  $B$  and  $\hat{B}$  than existing metrics [14]–[16]. Moreover, ACR provides a visual tool that facilitates error analysis. For instance, Fig. 1c reveals the MLS behavior of the HMM tracker.

We report experiments on datasets of Western classical, jazz and rock music, showing that ACR leads to new insights and a clearer picture of the performance of beat trackers than existing metrics. For reproducibility, we open source our code at <https://github.com/SunnyCYC/acr4mls>.

## II. EXISTING EVALUATION METRICS

In this section, we review the commonly used metrics.

### A. F1-Score and Continuity-based Evaluation Metrics

Given a reference beat sequence  $B$  and an estimated one  $\hat{B}$ , the F1-score [7], [27]–[29] considers an estimated beat  $\hat{b}_j \in \hat{B}$  as correct, if there is a reference beat  $b_i \in B$  that falls within a tolerance window  $\varepsilon$  (e.g.,  $\pm 70$ ms) to  $\hat{b}_j$ :  $|b_i - \hat{b}_j| \leq \varepsilon$ . Specifically, for a test set, we first calculate the precision  $P$  (the proportion of estimated beats that are considered correct) and the recall  $R$  (the proportion of reference beats that are found), and then compute  $F1 = 2PR / (P + R)$ .

Instead of evaluating each beat individually, *continuity*-based metrics [14]–[16], [25] further take into account the inter-beat interval (IBI), i.e., the temporal difference between two adjacent beats. For  $\hat{b}_j$  to be considered as correct, not only both  $\hat{b}_j$  and  $\hat{b}_{j-1}$  have to fall within the tolerance window around their corresponding reference beats (i.e.,  $b_i, b_{i-1}$ ), but also their IBIs have to be similar up to a factor of  $\gamma$  (i.e.,  $|(b_i - b_{i-1}) - (\hat{b}_j - \hat{b}_{j-1})| \leq \gamma(b_i - b_{i-1})$ ). Commonly, in the literature, the factor  $\gamma = 0.175$  is used. There are two main variants [16]: CMLt (correct metrical level) only compares  $B$  with  $\hat{B}$ , while AMLt (allowed metrical levels) synthetically creates variants of  $B$  that tap at a different yet plausible metrical level (e.g., half offbeat), compares each of the variant with  $\hat{B}$ , and outputs the best score. See Section II-C for details.

### B. L-Correct Detection

Grosche *et al.* [30] proposed a context-sensitive evaluation method, called “L-correct detection,” where the parameter  $L \in \mathbb{N}_{\geq 2}$  specifies the length of the temporal context in beats. The idea is to consider every reference beat  $b_i$  as an *instance* and check *instance by instance* (e.g., from  $i = 1$  onwards) if the subsequence  $B_{i,L} = \{b_i, b_{i+1}, \dots, b_{i+L-1}\}$  containing  $L$  consecutive reference beats starting from  $i$  can be fully matched by some subsequence of the estimated beats  $\hat{B}_{j,L} = \{\hat{b}_j, \hat{b}_{j+1}, \dots, \hat{b}_{j+L-1}\}$  that starts from  $j$ . This implies  $|b_{i+\tau} - \hat{b}_{j+\tau}| \leq \varepsilon$  for  $\tau \in \{0, 1, \dots, L-1\}$ , not allowing any false-positives or false-negatives. In case of a match, all the beats in  $B_{i,L}$  are considered *L-correct detected*. This ensures each L-correct detected reference beat to be within the group of at least  $L$  consecutive estimated beats. Using the number of *L-correct detected* beats, Grosche *et al.* [30] defined the L-correct recall, precision, and the F-measure  $F^L$ .

### C. Creation and Usage of Variants of Reference Beats

As shown in Table I, both AMLt and L-correct consider further the *half offbeat* variants of  $B$ , taking the middle of every two adjacent beats from  $B$  to create  $\beta_i \equiv (b_i + b_{i+1})/2$ . AMLt uses a new sequence  $B^{\text{off-1/2}} = (\beta_1, \beta_2, \dots)$  to compare with the whole sequence of  $\hat{B}$ , hence allowing a beat tracker to tap at half offbeats. However, according to whether  $\hat{B}$  matches better to  $B$  or to  $B^{\text{off-1/2}}$ , AMLt reports only the higher matching score. Namely, it assumes the beat tracker to tap either at onbeats, or at half offbeats, *throughout the piece*, thereby not allowing for MLS. In contrast, for the L-correct metric in [30], the evaluation is processed instance by instance. In other words, a beat  $b_i$  is considered L-correct if there is a matching estimated sequence either for  $B_{i,L}$  (onbeat case) or for  $B_{i,L}^{\text{off-1/2}}$  (offbeat case). This great flexibility is desirable to address the MLS issue.

AMLt and L-correct have their own strengths and weaknesses. AMLt cannot deal with MLS but it considers actually not only the offbeat but also the *subharmonic* and *harmonic* cases (i.e., integer divisions or multiples of a certain tempo, such as half tempo and double tempo [31]; see Section III for details). Furthermore, the L-correct metric partially addresses the MLS between onbeats and offbeats, but not the other metrical levels. The proposed ACR combines the advantages of AMLt and the L-correct metric by considering much more metric-level conditions while evaluating the matching of beats instance by instance.

## III. ANNOTATION COVERAGE RATIO

### A. Subharmonic and Harmonic Reference Beats

We describe below how we extend the L-correct metric to account for both tempo subharmonics and harmonics. Figure 2 visualizes the variants of reference beats adopted by ACR.

To account for of subharmonic tempi, we *sub-sample*  $B$  from the  $i$ -th instance onwards and derive the following:

$$\begin{aligned}
 B_{i,L}^{\text{half}} &= \{b_i, b_{i+2}, \dots, b_{i+2(L-1)}\}, \\
 B_{i,L}^{\text{third}} &= \{b_i, b_{i+3}, \dots, b_{i+3(L-1)}\}, \\
 B_{i,L}^{\text{quarter}} &= \{b_i, b_{i+4}, \dots, b_{i+4(L-1)}\}.
 \end{aligned} \tag{1}$$

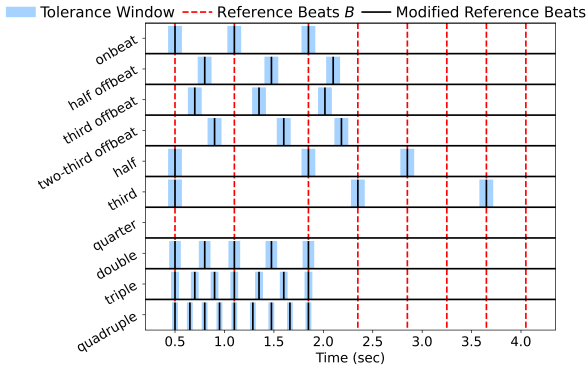


Fig. 2. Given a sequence of reference beats  $\{b_1, b_2, \dots, b_8\}$ , the figure shows the variants of the subsequence with length  $L = 3$  created for the first instance for different metrical levels. The grid corresponds to the reference beats.

If the sequence of reference beats  $B$  is not long enough to generate a modified reference beats of a specific metrical level ‘\*’, we assign  $B_{i,L}^* = \emptyset$  (see Figure 2 for  $*$  = quarter).

To account for harmonic tempi, we *upsample*  $B$  by linear interpolation, namely by taking the bisection points (for double tempo), trisection points (for triple tempo), and quadrisection points (for quadruple tempo) for the interval  $\overline{b_m b_{m+1}}$ ,  $m \in [i : i + L - 2]$  as additional reference beat positions. Due to the interpolation, the length of the resulting variants will be  $L' = L + (h - 1)(L - 1)$ , with  $h = 2, 3, 4$  corresponding to the double, triple, and quadruple cases. We therefore consider a longer context  $L'$  for L-correct detection of the harmonic cases. Denoting the  $\tau$ -th element of the variant created at instance  $i$  as  $B_{i,L'}^*[\tau]$ , where  $\tau$  starts from 0, we say that  $B_{i,L'}^*$  is L-correct detected if there exists an index  $j$  such that  $|B_{i,L'}^*[\tau] - \hat{b}_{j+\tau}| \leq \varepsilon$ , for  $\tau \in \{0, 1, \dots, L' - 1\}$ .

Conventionally, the tolerance window  $\varepsilon$  is usually set to a fixed value, such as 70ms. However, due to the upsampling, the IBI for some adjacent beats in  $B_{i,L'}^*$  (i.e.,  $B_{i,L'}^*[\tau + 1] - B_{i,L'}^*[\tau]$ ) might be already smaller than 70ms. To accommodate this, we instead compute the average IBI of  $B_{i,L'}^*$ , denoted as  $\overline{\text{IBI}}_{i,L'}^*$ , and use an adaptive  $\varepsilon = \min(70\text{ms}, \gamma \overline{\text{IBI}}_{i,L'}^*)$ , with  $\gamma = 0.175$  following the convention of CMLt/AMLt. Figure 2 shows how  $\varepsilon$  gets smaller for the tempo harmonic cases.

For highly expressive music, the additional beat positions created by linear interpolation may not be musically optimal [32]. We consider our approach mainly as an analysis tool and leave the discussion of better interpolation for future work.

### B. Visualization and Annotation Coverage

For every reference beat position  $b_i \in B$ , we propose to check if it is *L-correct detected* w.r.t. any of the ten metric-level conditions shown in Figure 2, and use a separate Boolean value to indicate the result. This way, we can plot the Boolean values for all the conditions in one graph to visualize whether and how each individual beat  $b_i$  is L-correct w.r.t. different conditions. As exemplified in Fig. 1, as long as  $b_i$  is L-correct detected w.r.t. a condition (e.g., offbeat), it will be ‘‘covered’’ by a colored bar in that row. We define the ‘‘annotation coverage ratio’’ (ACR) of each condition as the percentage of reference beats that are covered w.r.t. that condition.

TABLE II  
STATISTICS OF THE DATASETS.

Data	# tracks	Total duration	% stable tempi	Mean track tempo (BPM)
Maz-5 [33]	301	12h 27m	13.1%	125.39
RWC-Jazz [34]	50	3h 42m	74.8%	89.67
Rock [35]	200	12h 53m	81.4%	115.68

To have a more general view of the overall performance, we additionally define an *any tempo* case (shown in the bottom of Fig. 1b and 1c) that takes the union of the Boolean values of all metric-level cases per beat  $b_i$ , leading to ‘ACR-any.’ Similarly, we take the union of the three offbeat cases to derive ‘ACR-offbeat’ (shown in the second row of Fig. 1b and 1c). Note that, by definition, ACR cannot exceed 1.0 under any condition. To monitor and prevent inappropriately high MLS frequency, we define MLS ratio (MLSR) as the proportion of covered reference beats which switch metrical levels.<sup>2</sup>

## IV. EXPERIMENT SETUP

In our experiment, we use the three datasets listed in Table II to investigate the effectiveness of ACR for three important music genres—Western classical, jazz, and rock. Maz-5 [33] is a private collection of 301 recordings of five Chopin Mazurkas collected within the Mazurka Project [36] and annotated by Sapp [37]. RWC-Jazz (or ‘Jazz’ for short) contains 50 jazz recordings produced for the RWC Database [34]. The Rock dataset [35] contains 200 of the ‘‘500 Greatest Songs of All Time’’ listed the Rolling Stone magazine. Along with other statistics, Table II also shows the average tempo per track and the ‘‘percentage of stable tempi’’ [38] of the datasets.<sup>3</sup> We see that Maz-5 has the lowest number of stable tempi, while Jazz has the slowest average tempo.

As for the beat trackers, we adopt the two-block architecture mentioned in Section I and implement three approaches. They all use the popular open-source library, `madmom` [4], [13] to generate the beat activation functions  $\Delta(n)$ , but use different PPTs for converting  $\Delta(n)$  into the estimated beat positions  $\hat{B}$ . Specifically, we use a simple peak picking (SPPK) method [40], a DP-based tracker [3], [11], [41], and an HMM-based tracker [12], [13]. For SPPK, we use the ‘‘find\_peaks’’ function of the `scipy` library [40]. For HMM, we use the implementation in `madmom`, using the default setting. The DP requires a global tempo value to be provided beforehand. Following [30], we compute the global tempo from the mean IBI of the reference beats. Note that no training is involved in our experiment because the network provided by `madmom` is already trained and the three PPTs do not need any training.

## V. EXPERIMENTS

Table III shows the evaluation results of the beat trackers using the proposed ACR and existing metrics including F1-

<sup>2</sup>One may also increase  $L$  value to achieve stricter metric-level criteria.  
<sup>3</sup>Following [38], after converting all IBIs of each dataset into tempo values and dividing these values with their corresponding mean track tempo, we can derive the normalized tempi and calculate the ‘‘percentage of stable tempi’’ as the proportion of normalized tempi that falls in the commonly adopted  $\pm 4\%$  tolerance interval of stable tempi [39].

TABLE III  
THE PERFORMANCE OF DIFFERENT POST-PROCESSING TRACKERS (PPTs) ON DIFFERENT DATASETS USING DIFFERENT EVALUATION METRICS.

Data	PPT	F1	R	P	CMLt	AMLt	L-correct $F^L(L=2)$	ACR (L=2)									MLSR
								any	onbeat	offbeat	subharmonics			harmonics			
											half	third	quarter	double	triple	quadruple	
Maz-5	DP	0.488	0.493	0.468	0.315	0.315	0.590	0.598	0.383	0.181	0.025	0.001	0.000	0.013	0.002	0.001	0.006
	HMM	0.499	0.393	0.752	0.098	0.248	0.128	0.726	0.112	0.007	0.283	0.235	0.128	0.003	0.001	0.000	0.002
	SPPK	0.822	0.754	0.918	0.569	0.570	0.743	0.904	0.682	0.003	0.276	0.070	0.025	0.018	0.000	0.000	0.005
Jazz	DP	0.781	0.796	0.765	0.787	0.845	0.895	0.915	0.790	0.124	0.000	0.000	0.000	0.001	0.000	0.000	0.000
	HMM	0.797	0.889	0.746	0.649	0.803	0.674	0.903	0.669	0.026	0.000	0.000	0.000	0.192	0.000	0.016	0.000
	SPPK	0.680	0.822	0.601	0.336	0.568	0.434	0.774	0.452	0.010	0.029	0.004	0.000	0.354	0.002	0.017	0.001
Rock	DP	0.961	0.974	0.949	0.947	0.948	0.968	0.982	0.970	0.011	0.000	0.000	0.000	0.000	0.000	0.000	0.000
	HMM	0.948	0.982	0.925	0.886	0.925	0.901	0.977	0.914	0.002	0.005	0.000	0.000	0.047	0.010	0.000	0.000
	SPPK	0.915	0.972	0.879	0.771	0.827	0.820	0.962	0.843	0.001	0.009	0.002	0.001	0.129	0.025	0.000	0.003

score, CMLt/AMLt, and L-correct ( $F^L$ ). According to the result of the existing metrics, we can see that beat tracking seems to be easier for music with steady tempo (cf. Table II): for example, all the three trackers have F1 scores above 0.900 on the Rock dataset. However, a closer look reveals that there is inconsistency among the results of the existing metrics. We make the following observations. **O1**: For both Maz-5 and Jazz, HMM slightly outperforms DP in F1, but, without explanation, DP gets much better result in L-correct and CMLt/AMLt metrics. **O2**: From the low  $F^L$  (0.128), low recall (0.393), and high precision (0.752) of HMM for Maz-5, we conjecture that HMM may tend to tap at slower tempi than the reference beats (therefore cannot match the reference beats consecutively for  $L = 2$ ), but none of the existing metrics clearly support this. **O3**: Similarly, compared to DP, HMM has lower L-correct, higher recall and lower precision for Jazz, suggesting that HMM may tap at faster tempi here, but again more evidence is needed.

Our newly proposed ACR metrics provide explanations to these observations and lead to new insights. For example, from the ACR scores for Maz-5 in Table III, one can understand that HMM not only taps at slower tempi but also switches among different tempo subharmonics, which explains its extremely low scores w.r.t. L-correct and CMLt/AMLt metrics. Similarly, the ACR scores show that HMM does tap faster and mainly at double tempo for Jazz. It seems that HMM tends to tap slower at subharmonics for faster tracks (e.g., Maz-5) and tap faster at harmonics for slower tracks (e.g., Jazz), explaining its inferior scores in existing metrics such as L-correct compared to DP. Besides, from the ACR of DP, we see that although DP is exempt from the harmonic/subharmonic “errors,” it instead suffers from offbeat “errors,” as also reported in [30]. In sum, for challenging music pieces, existing beat trackers have not learned to adapt to the tempo changes, and could produce different types of metric-level or phase-related “errors” that lead to poor scores not fully explainable by existing metrics. The ACR reveals these metric-level behaviors and provides additional evaluation perspectives. For example, according to ACR-any, if MLS is allowed in the middle of a music piece, HMM is actually not inferior to DP.

Table IV further shows the ACR results for the metrical level of onbeat and ‘any tempo’ with  $L = 2, 3, 4$ , which provides

TABLE IV  
CONTEXT SENSITIVE EVALUATION RESULT WITH  $L = 2, 3, 4$ .

Data	PPT	ACR-onbeat			ACR-any		
		L=2	L=3	L=4	L=2	L=3	L=4
Maz-5	DP	0.383	0.310	0.243	0.598	0.414	0.294
	HMM	0.112	0.102	0.096	0.726	0.597	0.484
	SPPK	0.682	0.622	0.566	0.904	0.721	0.622
Jazz	DP	0.790	0.787	0.784	0.915	0.904	0.895
	HMM	0.669	0.667	0.666	0.903	0.893	0.887
	SPPK	0.452	0.366	0.324	0.774	0.658	0.576
Rock	DP	0.970	0.969	0.969	0.982	0.979	0.977
	HMM	0.914	0.914	0.914	0.977	0.977	0.976
	SPPK	0.843	0.813	0.797	0.962	0.931	0.906

insights for how well a beat tracker can handle the temporal context of several consecutive beats, rather than looking at the beats individually. For Maz-5, using either ACR-onbeat or ACR-any, the scores for the three PPTs all drop remarkably as  $L$  increases, indicating that none of them can handle well the temporal context. On the contrary, for Jazz and Rock, the scores of DP and HMM remain almost unchanged as  $L$  increases, indicating that most of the estimated beats are predicted correctly for a sequence of consecutive beats rather than individually. From the remarkably higher ACR-any values compared to ACR-onbeat values of DP and HMM for Jazz, we also see that they can handle more than 90% of the temporal context well (i.e., ACR-any > 0.90) but with metric-level switching. From the performance degradation of SPPK as  $L$  increases, we see that SPPK, which makes estimation merely based on individual activation peaks, cannot handle temporal context as HMM and DP can do for Jazz and Rock.

## VI. CONCLUSION

In this paper, we introduced an analysis method, called annotation coverage ratio (ACR), to reveal the metric-level switching behaviors of beat trackers. With experiments on datasets of three major music genres, we demonstrate how the proposed ACR can compensate for the inadequacy of existing evaluation metrics and facilitate further analysis of beat tracking for challenging genres of music (e.g., classical and jazz). We hope this work can contribute to the development of more advanced and general beat trackers.

## REFERENCES

- [1] M. Goto and Y. Muraoka, "A beat tracking system for acoustic signals of music," in *Proc. ACM Int. Conf. Multimedia*, 1994, pp. 365–372.
- [2] S. Dixon and E. Cambouropoulos, "Beat tracking with musical knowledge," in *Proc. Eur. Conf. on Artificial Intelligence*, 2000, p. 626–630.
- [3] D. P. Ellis, "Beat tracking by dynamic programming," *J. New Music Res.*, vol. 36, no. 1, pp. 51–60, 2007.
- [4] S. Böck, F. Krebs, and G. Widmer, "Joint beat and downbeat tracking with recurrent neural networks," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2016, pp. 255–261.
- [5] M. Fuentes, B. McFee, H. C. Crayencour, S. Essid, and J. P. Bello, "Analysis of common design choices in deep learning systems for downbeat tracking," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2018, pp. 106–112.
- [6] M. E. P. Davies and S. Böck, "Temporal convolutional networks for musical audio beat tracking," in *Proc. Eur. Signal Process. Conf.*, 2019, pp. 1–5.
- [7] S. Böck and M. E. P. Davies, "Deconstruct, analyse, reconstruct: How to improve tempo, beat, and downbeat estimation," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2020, p. 574–582.
- [8] M. Heydari and Z. Duan, "Don't look back: An online beat tracking method using RNN and enhanced particle filtering," in *IEEE Int. Conf. Acoust. Speech Signal Process.*, 2021, pp. 236–240.
- [9] C.-Y. Chiu, A. W.-Y. Su, and Y.-H. Yang, "Drum-aware ensemble architecture for improved joint musical beat and downbeat tracking," *IEEE Signal Processing Letters*, pp. 1100–1104, 2021.
- [10] Y.-N. Hung, J.-C. Wang, X. Song, W.-T. Lu, and M. Won, "Modeling beats and downbeats with a time-frequency Transformer," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 2022, pp. 401–405.
- [11] M. Müller, *Fundamentals of Music Processing – Using Python and Jupyter Notebooks*, 2nd ed. Springer Verlag, 2021.
- [12] F. Krebs, S. Böck, and G. Widmer, "An efficient state-space model for joint tempo and meter tracking," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2015, pp. 72–78.
- [13] S. Böck, F. Korzeniowski, J. Schlüter, F. Krebs, and G. Widmer, "Madmom: A new Python audio and music signal processing library," in *Proc. ACM Multimed. Conf.*, 2016, pp. 1174–1178.
- [14] M. Macleod and S. Hainsworth, "Particle filtering applied to musical tempo tracking," *EURASIP J. Advances in Signal Processing*, vol. 2004, no. 15, pp. 2385–2395, 2004.
- [15] A. Klapuri, A. Eronen, and J. Astola, "Analysis of the meter of acoustic musical signals," *IEEE Trans. Audio, Speech, and Language Process.*, vol. 14, pp. 342 – 355, 2006.
- [16] M. E. P. Davies, N. D. Quintela, and M. Plumbley, "Evaluation methods for musical audio beat tracking algorithms," in *Queen Mary University of London, Centre for Digital Music, Tech. Rep. CADM-TR-09-06*, 2009.
- [17] M. McKinney and D. Moelants, "Ambiguity in tempo perception: What draws listeners to different metrical levels?" *Music Perception*, vol. 24, pp. 155–166, 2006.
- [18] M. McKinney, D. Moelants, M. E. P. Davies, and A. Klapuri, "Evaluation of audio beat tracking and music tempo extraction algorithms," *J. New Music Research*, vol. 36, pp. 1–16, 2007.
- [19] E. Cano, F. Ángel, L. Gil, J. Zapata, A. Escamilla, J. Alzate, and M. Betancur, "Sesquialtera in the colombian bambuco: Perception and estimation of beat and meter," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2020.
- [20] E. Cano, F. Mora-Ángel, L. Gil, J. Zapata, A. Escamilla, J. Alzate, and M. Betancur, "Sesquialtera in the colombian bambuco: Perception and estimation of beat and meter – extended version," *Trans. of the Int. Soc. Music Inf. Retr.*, vol. 4, p. 248, 2021.
- [21] H. Schreiber, J. Urbano, and M. Müller, "Global music tempo estimation: Are we done yet?" *Transactions of the International Society for Music Information Retrieval (TISMIR)*, vol. 3, no. 1, pp. 111–125, 2020.
- [22] M. A. Miguel and D. F. Slezak, "Modeling beat uncertainty as a 2D distribution of period and phase: A MIR task proposal," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2021, pp. 452–459.
- [23] M. E. P. Davies and M. Plumbley, "Context-dependent beat tracking of musical audio," *IEEE Trans. Audio, Speech, and Language Process.*, vol. 15, pp. 1009 – 1020, 2007.
- [24] A. S. Pinto, S. Böck, J. S. Cardoso, and M. E. P. Davies, "User-driven fine-tuning for beat tracking," *Electronics*, vol. 10, no. 13, 2021.
- [25] M. E. P. Davies, S. Böck, and M. Fuentes, *Tempo, beat and downbeat estimation*. Proc. Int. Soc. Music Inf. Retr. Conf., 2021. [Online]. Available: <https://tempobeatdownbeat.github.io/tutorial/intro.html>
- [26] C.-Y. Chiu, M. Müller, M. E. P. Davies, A. W.-Y. Su, and Y.-H. Yang, "Local temporal expectation-based beat tracking for expressive classical music," *IEEE/ACM Trans. Audio, Speech, and Language Process.*, 2022, under review.
- [27] M. E. P. Davies and S. Böck, "Evaluating the evaluation measures for beat tracking," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2014.
- [28] B. Jia, J. Lv, and D. Liu, "Deep learning-based automatic downbeat tracking: a brief review," *Multimedia Systems*, vol. 25, no. 6, pp. 617–638, 2019.
- [29] T. Oyama, R. Ishizuka, and K. Yoshii, "Phase-aware joint beat and downbeat estimation based on periodicity of metrical structure," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2021, pp. 493–499.
- [30] P. Grosche and M. Müller, "Extracting predominant local pulse information from music recordings," *IEEE Trans. Audio, Speech Lang. Process.*, vol. 19, no. 6, pp. 1688–1701, 2011.
- [31] M. Müller and V. Arifi-Müller. Tempo and beat. [Online]. Available: [https://www.audiolabs-erlangen.de/resources/MIR/FMP/C6/C6S2\\_TempoBeat.html](https://www.audiolabs-erlangen.de/resources/MIR/FMP/C6/C6S2_TempoBeat.html)
- [32] P. Desain and H. Honing, "Tempo curves considered harmful," *Contemporary Music Review*, vol. 7, pp. 123–138, 1993.
- [33] P. Grosche, M. Müller, and C. S. Sapp, "What makes beat tracking difficult? A case study on Chopin Mazurkas," *Proc. Int. Soc. Music Inf. Retr. Conf.*, no. January, pp. 649–654, 2010.
- [34] M. Goto, H. Hashiguchi, T. Nishimura, and R. Oka, "RWC Music Database: Popular, Classical, and Jazz music databases," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2002, pp. 287–288.
- [35] T. de Clercq and D. Temperley, "A corpus analysis of rock harmony," *Popular Music*, vol. 30, no. 1, p. 47–70, 2011.
- [36] (2010) The Mazurka Project. [Online]. Available: <http://mazurka.org.uk/>
- [37] C. Sapp, "Hybrid numeric/frank similarity metrics for musical performance analysis," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2008, pp. 501–506.
- [38] H. Schreiber, F. Zalkow, and M. Müller, "Modeling and estimating local tempo: A case study on Chopin's Mazurkas," in *Proc. Int. Soc. Music Inf. Retr. Conf.*, 2020, pp. 773–779.
- [39] F. Gouyon, A. Klapuri, S. Dixon, M. Alonso, G. Tzanetakis, C. Uhle, and P. Cano, "An experimental comparison of audio tempo induction algorithms," *IEEE Trans. Audio, Speech, and Language Process.*, vol. 14, no. 5, pp. 1832–1844, 2006.
- [40] P. Virtanen *et al.*, "SciPy 1.0: Fundamental algorithms for scientific computing in Python," *Nature Methods*, vol. 17, pp. 261–272, 2020.
- [41] M. Müller and F. Zalkow, "libfmp: A Python package for fundamentals of music processing," *J. Open Source Software*, vol. 6, no. 63, p. 3326, 2021.