

TOWARDS HUMAN-IN-THE-LOOP ONSET DETECTION: A TRANSFER LEARNING APPROACH FOR *MARACATU*

António Sá Pinto

Faculdade de Engenharia da Universidade do Porto, Porto, Portugal

INESC TEC, Porto, Portugal

asapinto@fe.up.pt

ABSTRACT

We explore transfer learning strategies for musical onset detection in the Afro-Brazilian *Maracatu* tradition, which features complex rhythmic patterns that challenge conventional models. We adapt two Temporal Convolutional Network architectures: one pre-trained for onset detection (intra-task) and another for beat tracking (inter-task). Using only 5-second annotated snippets per instrument, we fine-tune these models through layer-wise retraining strategies for five traditional percussion instruments. Our results demonstrate significant improvements over baseline performance, with F1 scores reaching up to 0.998 in the intra-task setting and improvements of over 50 percentage points in best-case scenarios. The cross-task adaptation proves particularly effective for time-keeping instruments, where onsets naturally align with beat positions. The optimal fine-tuning configuration varies by instrument, highlighting the importance of instrument-specific adaptation strategies. This approach addresses the challenges of underrepresented musical traditions, offering an efficient human-in-the-loop methodology that minimizes annotation effort while maximizing performance. Our findings contribute to more inclusive music information retrieval tools applicable beyond Western musical contexts.

1. INTRODUCTION

Accurately identifying the precise moment when a musical note begins remains one of the fundamental challenges in audio signal processing. This task, known as musical onset detection, serves as a cornerstone for numerous Music Information Retrieval (MIR) applications. Onset detection has historically been essential for rhythmic analysis, notably in beat tracking systems [1–3]. While end-to-end learning models have recently bypassed this explicit step in some contexts, onset detection continues to be critical for diverse applications such as score following [4], music segmentation [5], and polyphonic music transcription [6].

The methodological evolution of onset detection mirrors broader trends in MIR research. Early approaches re-

lied on signal processing techniques to identify significant changes in audio properties [7, 8], followed by the introduction of feature-based machine learning methods [9, 10]. The field then shifted toward neural network architectures, beginning with Recurrent Neural Networks (RNNs) [11] and advancing to Convolutional Neural Networks (CNNs) [12], which extract relevant features directly from raw audio or spectral representations. Despite impressive advances in performance metrics (with top models achieving F1 scores approaching 90% in recent evaluations¹), significant challenges persist in onset detection. In particular, accurately detecting soft onsets remains difficult even for advanced models [13]. Moreover, these data-driven approaches introduce additional challenges related to training data requirements and generalizability.

The effectiveness of supervised learning models hinges on the quality and diversity of training data [14]. Current systems experience performance drops when analysing non-Western musical traditions or rare instruments, primarily due to insufficient representation in existing datasets. Addressing these gaps requires costly annotation efforts that demand both domain-specific and culturally-informed expertise [15], further complicating dataset curation. Furthermore, the annotation process itself reveals limitations: manual labelling of onsets is prone to human error and inconsistencies [16], with even isolated percussive signals proving difficult to label precisely [17]. These constraints restrict the practical deployment of state-of-the-art systems in diverse musical contexts, pointing to the need for more adaptable strategies.

Moving beyond the specific challenges of onset detection, MIR research has employed several adaptive strategies within rhythm analysis tasks. Informed methods leverage a priori knowledge about rhythmic content for tasks such as beat tracking [18] and metre determination [19], which, while effective in specific genres, lack generalizability. Transfer learning leverages knowledge across domains, with examples including adaptations of mainstream beat-tracking models to Greek folk music [20] and facilitating adaptive rhythm microtiming generation [21]. Additionally, user-centric approaches like Active Learning and Few-Shot Learning optimize learning through strategic sample selection, enhancing adaptability in polyphonic drum transcription [22, 23] and enabling in-



¹ MIREX 2018, at https://nema.lis.illinois.edu/nema_out/mirex2018/results/aod/

teractive refinement for onset detection [24] and beat tracking [25]. This shift toward user involvement exemplifies the current human-centred landscape of MIR, recognizing users’ essential role in data-driven systems [26]. The integration of human expertise into computational frameworks provides a promising avenue when existing solutions prove insufficient.

Recent research has explored incorporating user-provided information to enhance beat tracking performance. Techniques such as high-level model parameterization [27] and integrating user-annotated data snippets in a fine-tuning cycle [28] have shown promise for improving state-of-the-art accuracy. These methods are particularly effective in addressing challenges in underrepresented musical contexts, where conventional MIR techniques underperform. Such approaches have proven instrumental in the creation of the *Maracatu* onset dataset [17], metre determination in Latin-American music [29], and beat tracking in highly challenging music signals [30]. The implementation of transfer learning for these tasks varies considerably: while some approaches retrain only final layers to leverage basic rhythmic representations [20, 31], others target input and output layers for instrument-specific adaptation [17], and some retrain entire networks [28]. Despite these varied strategies, no studies have empirically evaluated the impact of layer-wise retraining on model performance, leaving this critical question unexplored.

Building on this foundation, this paper explores a user-driven transfer learning approach for onset detection, focusing on the Afro-Brazilian tradition of *Maracatu*. We use the eponymous dataset [17], which features complex rhythms and unique instrumental acoustic characteristics that cause leading models to struggle with achieving satisfactory performance.

Our methodology involves adapting a deep neural network for each instrument in the “terno”, the percussion ensemble central to *Maracatu*’s rhythm, based on a short annotated snippet per instrument. Through these instrument-specific adaptations, we demonstrate an effective and straightforward method to enhance state-of-the-art performance. We investigate two distinct transfer learning scenarios: one with a model initially trained for onset detection, and another novel approach adapting a beat-tracking model to onset detection. This extends previous research [17, 32] by exploring cross-task feature transferability and leveraging more complex models trained on larger datasets. Furthermore, we systematically evaluate layer-wise retraining strategies, examining the effectiveness of freezing different layer groups to identify optimal configurations for *Maracatu* onset detection.

2. METHODOLOGY

Our approach addresses the limitations of existing models in non-mainstream signals by integrating user-provided short annotated snippets. We adapt the human-in-the-loop method proposed by Pinto et al. for beat tracking [27, 28, 30] to the task of onset detection, leveraging state-of-the-art models through *in-situ* fine-tuning. This

user-centred methodology eliminates the need for extensive training from scratch, enabling end-users to swiftly obtain high-quality onset estimates that align with their judgments.

For onset detection in monotonimbral signals, we adapt neural networks to each instrument’s unique acoustic characteristics using just a single 5-second annotated snippet per instrument as the fine-tuning target. This approach demonstrates both minimal annotation effort and rapid adaptation cycles, yielding instrument-specific networks optimized for their corresponding acoustic properties while remaining computationally feasible for standard resources. While our method is applicable to various DNN architectures, this study employs Temporal Convolutional Network (TCN)-based models for their efficient retraining capabilities. The TCN’s performance in onset detection tasks is comparable to state-of-the-art models, as demonstrated in Section 3.1, making it suitable for our investigation.

We explore two transfer learning scenarios: an intra-task setting using a TCN onset detection model [32] and an inter-task setting that adapts a TCN beat tracking model [33] to onset detection. This inter-task approach can be framed as a domain adaptation problem, where a model trained for beat tracking is repurposed for onset detection. Given the inherent relationship between beats and onsets, this adaptation may benefit from the typically broader training data available for beat tracking models. To the best of our knowledge, this is the first study to explore domain adaptation from beat tracking to onset detection.

Furthermore, onset detection’s unambiguous objective, when contrasted with the multifaceted nature of beat tracking, allows for clearer adaptation targets and, consequently, more straightforward interpretation of results. This motivated us to extend previous research by examining layer-wise retraining strategies. We systematically freeze different segments of the 15-layer TCN architectures, from the initial convolutional layers with small receptive fields to the deeper layers with larger dilation rates and wider receptive fields. In total, our experimental cycle comprises 150 fine-tuning cycles (15 layer configurations \times 5 instruments \times 2 models). Through this comprehensive evaluation, we aim to investigate feature transferability between related rhythm analysis tasks and systematically assess the impact of different layer freezing configurations.

In line with open science principles [34], we provide a GitHub repository with our code and detailed results, including per-file evaluation metrics for all configurations and higher-resolution figures for detailed analysis². The remainder of this section outlines the *Maracatu* dataset composition, experimental settings, base models’ description, and fine-tuning and evaluation details.

² <https://github.com/asapsmc/HIILOnsetDetection>

2.1 Dataset

*Maracatu de baque solto*³, also known as *Maracatu* “rural”, is a vibrant carnival performance from Pernambuco, Northeast Brazil, combining music, poetry, and dance [36]. The rhythmic nucleus of *Maracatu*, known as the “terno” ensemble, consists of five percussionists playing traditional handmade instruments: *cuica*, *gonge-lo*, *tarol*, *mineiro*, and *tambor-hi*. The *Maracatu* dataset [17] captures these instruments using contact microphones for largely isolated per-instrument tracks, recorded during a fixed location performance and comprising 34 individual pieces totalling approximately 33 minutes⁴.

Maracatu features two main rhythmic patterns: “marcha” and “samba”, characterized by fast tempi of approximately 165 and 180 beats-per-minute (bpm), respectively. This rapid pace creates a complex timing profile across the ensemble. Time-keeping instruments (*cuica* and *gonge-lo*) maintain rhythmic stability despite their sporadic use, with a mean onset count of around 4,700 (2.5 annotations per second). In contrast, the “voicing” instruments (*tarol*, *mineiro*, and *tambor-hi*) play more expressive roles, resulting in a higher mean onset count of approximately 16,600 (8.9 annotations per second).

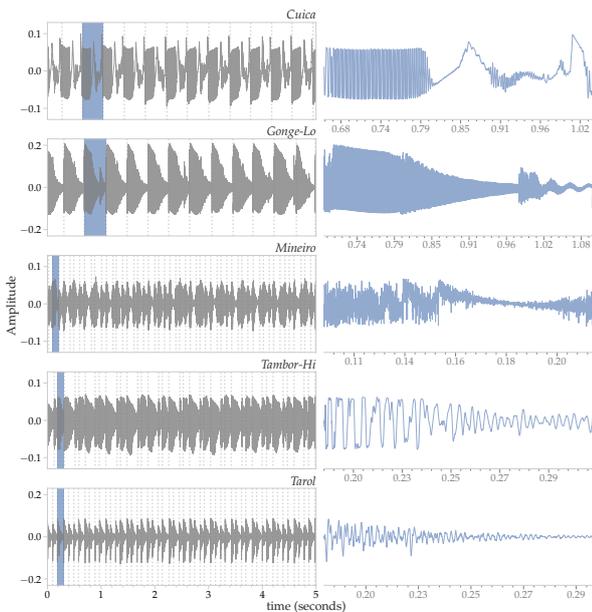


Figure 1. Onset-annotated waveforms for the *Maracatu* instruments. Left: 5-seconds fine-tuning snippet; Right: Zoomed in waveform, from the second onset to the sample before third onset (in blue).

The intricacy of these rhythms and distinct waveform shapes, as illustrated in Figure 1, complicates onset detection and annotation. The *mineiro* exemplifies this challenge with its unusual waveform characteristics, which led

³ Hereafter referred to as *Maracatu*, this genre should be distinguished from *Maracatu de baque virado* (or “Nação”). Both share African origins and certain musical similarities, but differ significantly in instrumentation, practice, and narrative [35].

⁴ While the original dataset contains 34 files per instrument, we excluded *Instrument_34* files across all sub-datasets due to a corrupted *Mineiro_34* file.

to its exclusion from microtiming analysis in the original dataset creation study due to annotation difficulties [17]. Combined with the under-representation of these instruments in available model training data, these factors create substantial obstacles for both human annotators and automated systems. The *Maracatu* dataset thus provides an ideal test bed for our human-in-the-loop strategy, extending the approach previously employed in the dataset’s creation.

2.2 Base Models

This study employs two pre-trained models, both derived from the TCN architecture proposed by Davies and Böck [37]. For the intra-task setting, we use a modified version of the original TCN model with an additional 11th dilation level [32], trained from scratch on the *OnsetDB* dataset [4] for onset detection. In the inter-task scenario, we utilize an adaptation of the [33] multitask network, modified by masking its tempo and downbeat loss to function as a single-task (beat) network, trained on various beat-tracking datasets. Hereafter, we refer to these models as TCN_{v1} and TCN_{v2} , respectively.

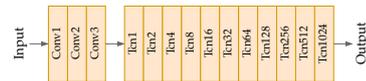


Figure 2. High-level architecture shared by the TCN_{v1} and TCN_{v2} models. Both follow the same layer sequence and depth, but differ in convolutional filter configuration, resulting in distinct receptive fields and overall model sizes.

As illustrated in Figure 2, both models share the same high-level architecture and signal conditioning stages, but their implementations differ significantly. TCN_{v1} consists of three convolutional layers with 16 filters and filter shapes of 3×3 , 3×3 , and 1×8 , with max pooling over three frequency bins after the first two layers. In contrast, TCN_{v2} employs three convolutional layers with 20 filters and filter shapes of 3×3 , 1×10 , and 3×3 , each followed by max pooling over three frequency bins. Both architectures use dropout after each convolutional stage. The ensuing TCN block operates non-casually and consists of 11 dilation levels, 16 filters, and a kernel size of 5. The TCN_{v1} model comprises 21,890 parameters, while the TCN_{v2} model has 116,302 parameters. The original training procedures also differed slightly in optimization techniques: TCN_{v1} employed a standard *Adam* optimizer, whereas TCN_{v2} used a *Rectified Adam* plus *Lookahead* approach.

2.3 Fine-tuning

For both intra-task and inter-task transfer learning settings, we adopt the fine-tuning strategy described in [28], using a 5-second annotated sample per instrument to demonstrate minimal annotation effort. Each base model is fine-tuned for 50 epochs with the learning rate reduced to one-quarter of the original value, maintaining the original optimizers for seamless training continuation. Early stopping

and learning rate reduction mechanisms were not implemented as these parameters proved sufficient for convergence given the short training duration and small dataset size. Given our systematic layer-wise analysis comprising 150 fine-tuning cycles, we omitted data augmentation and additional hyperparameter optimization to maintain experimental tractability and support isolated analysis of how layer-wise freezing strategies relate to each instrument’s acoustic characteristics.

We evaluate all possible fine-tuning configurations, denoted as ft_{A-B} , where A and B indicate the starting and ending layers of the frozen section, respectively. The output layer is always updated and thus excluded from this notation. We explore configurations from $ft_{Conv1\dots3}$ to $ft_{Tcn1\dots1024}$, including the fully trainable configuration ft . These are compared with the intra-task baseline bsl and the inter-task baseline bsl^* .

2.4 Evaluation

The network output is an onset activation function with a 10-millisecond (ms) temporal resolution. We apply the standard `madmom` peak-picking algorithm to obtain onset estimates. Performance is evaluated using the F1 metric with the default 25 ms tolerance window [4]. We implement a holdout validation approach where, for each instrument, we extract a 5-second segment from the first file (*Instrument_01*) for fine-tuning and then exclude this entire file from the evaluation set to prevent data leakage. This ensures unbiased assessment of the instrument-adapted models by evaluating performance on the remaining 32 files per instrument.

3. EXPERIMENTS AND RESULTS

3.1 Preliminary Model Analysis

To contextualize our approach, we first compare the performance of our base TCN models with previous state-of-the-art methods on the *OnsetDB* dataset [4]. Our base models, `TCNv1` and `TCNv2`, achieve F1 scores of 0.907 and 0.340, respectively. The lower performance of `TCNv2` is expected, as it was originally trained for beat tracking rather than onset detection. The `madmomRNN` and `madmomCNN` models, pre-trained and provided as inference-ready models in the `madmom` package [38], achieve F1 scores of 0.849 and 0.913, respectively. However, it is important to note that these evaluations were conducted without knowledge of the original training/test splits used for these pre-trained models, creating potential data leakage that may lead to an overestimation of their performance. The 2nd generation onset CNN [39] remains the established benchmark, with a reported F1 score of 0.903, verified through k-fold cross-validation. Unlike the `madmom` models, our TCN models were evaluated under the same validation conditions as the 2nd gen CNN, ensuring comparability. These results indicate that `TCNv1` is competitive with the current state of the art in onset detection.

Table 1. Representative configurations demonstrating improvements across transfer learning settings.

	<i>Onset-to-Onset</i>		<i>Beat-to-Onset</i>			
	Adapted (best)	bsl	Adapted (best)	bsl^*		
<i>Cuica</i>	ft_{Tcn16}	0.985	0.477	ft	0.955	0.429
<i>Gonge-Lo</i>	$ft_{Tcn2/4/16}$	0.998	0.508	ft	0.956	0.892
<i>Mineiro</i>	ft_{Tcn16}	0.972	0.946	ft_{Tcn8}	0.790	0.193
<i>Tambor-Hi</i>	$ft_{/Tcn1024}$	0.978	0.965	ft_{Tcn1}	0.723	0.443
<i>Tarol</i>	ft_{Conv3}	0.997	0.993	ft	0.884	0.139

3.2 Onset-to-Onset Transfer Learning Results

Figure 3 (top) presents the F1 scores obtained for each fine-tuning configuration in comparison to the baseline. The results can be grouped based on the rhythmic role of the instruments: time-keeping (*cuica* and *gonge-lo*) vs. voicing (*tarol*, *mineiro*, and *tambor-hi*).

For time-keeping instruments, the baseline performance is moderate (F1 \approx 0.5), but fine-tuning yields significant improvements, with scores reaching the 0.8–1.0 range. In contrast, expressive instruments exhibit higher initial F1 scores (\approx 0.9–1.0), which limits the relative improvement. This disparity can be attributed to the conventional nature of *tarol* and *tambor-hi*, which are more aligned with the training data, whereas *cuica* and *gonge-lo* diverge more in terms of acoustic characteristics. An exception is *mineiro*, which achieves a relatively high baseline score despite its distinct waveform characteristics. However, the reported lower precision of these ground-truth annotations [17] complicates direct performance comparisons.

Table 1 presents high-performing configurations to demonstrate the achievable improvements across instruments. The ft_{Tcn16} model achieves the highest accuracy for *cuica* and *mineiro* (0.985 and 0.972, respectively), while ft_{Tcn2} , ft_{Tcn4} , and ft_{Tcn16} all achieve the highest F1 score for *gonge-lo* (F1 = 0.998). For *tambor-hi*, the best performance is obtained with both $ft_{Tcn1024}$ and ft (F1 = 0.978). For *tarol*, the highest F1 score (0.997) is achieved with ft_{Conv3} , though many configurations show comparable performance with marginal differences. These configurations consistently outperform the baseline, with the most notable gains observed in *cuica* and *gonge-lo*, where F1 improvements exceed 50 percentage points (p.p.).

In summary, all instruments benefit from adaptation, as most fine-tuned configurations—and in particular, the best for each instrument—consistently outperform the baseline. The improvement is especially pronounced for time-keeping instruments (*cuica* and *gonge-lo*), likely due to their lower baseline accuracy, which allows more room for improvement, and the relative ease of detecting sparser onsets compared to those that are closely clustered in time, even though onset density remains well above the network’s temporal resolution of 10 ms. The optimal freeze configuration varies by instrument, with no clear global trend. However, some patterns emerge: for voicing instruments, full-network fine-tuning (ft) ranks among the top-performing configurations, whereas it degrades performance for time-keeping instruments.

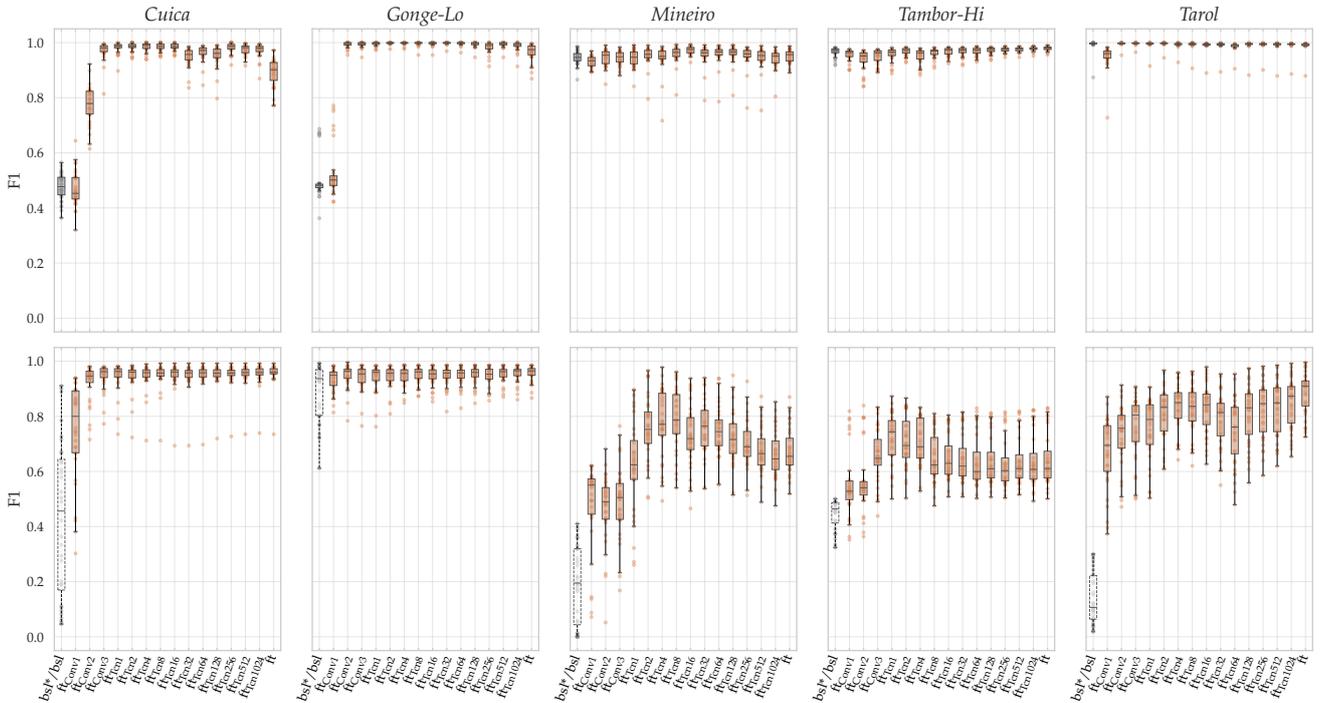


Figure 3. Distribution of F1 scores per layer-wise configuration under two transfer learning settings: *Onset-to-Onset* (top), where fine-tuned models are compared against their baseline, and *Beat-to-Onset* (bottom), where we assess cross-task versus within-task transfer learning, with comparable performance observed for time-keeping instruments.

3.3 *Beat-to-Onset* Transfer Learning Results

In this section, we focus on a domain adaptation, where a model pre-trained for beat tracking is adapted for onset detection. Unlike the previous setting, the goal here is not to compare fine-tuned models to their baseline, as this originates from a different task. We also refrain from an in-depth analysis of mean F1 scores across datasets, given their limited interpretative value. Instead, we assess whether models fine-tuned in this setting achieve results comparable to those in the onset-to-onset transfer learning scenario. Figure 3 (bottom) provides an overview of the results.

Time-keeping instruments, such as *cuica* and *gonge-lo*, achieve relatively high baseline ($bs1^*$) accuracies, likely due to the alignment between their onsets and beat locations. Adaptation improves accuracy across all instruments, confirming the feasibility of beat-to-onset transfer learning. However, while the fine-tuned models consistently outperform the beat-tracking baseline, direct comparisons to the onset-to-onset setting reveal performance disparities that vary by instrument. Specifically, for time-keeping instruments, performance remains nearly identical across both transfer learning scenarios, with differences of only 1.6 p.p. for *cuica* and 3.7 p.p. for *gonge-lo*. In contrast, voicing instruments exhibit progressively larger discrepancies, with F1-score differences of 11.3 p.p. for *tarol*, 27.5 p.p. for *mineiro*, and the largest gap of 32.2 p.p. for *tambor-hi*.

Closer inspection of the layer-wise results reveals additional patterns. The accuracy generally increases as more layers are fine-tuned up to the 3rd or 4th dilation level, be-

yond which no further gains are observed. However, this trend does not hold for *tarol*, where deeper fine-tuning leads to additional performance improvements. These observations highlight that, while fine-tuning is beneficial across all cases, the optimal retraining depth remains instrument-dependent.

Altogether, the results indicate that feature transferability from beat tracking to onset detection is more effective for time-keeping instruments than for voicing instruments. Specifically, *gonge-lo* exhibits a clearly higher baseline F1 accuracy in the beat-to-onset setting compared to its onset-to-onset counterpart (0.892 vs. 0.508), while *cuica* achieves a comparable performance (0.429 vs. 0.477), as reported in Table 1. This enhanced cross-task adaptability arises from the metrical function of time-keeping instruments: their onsets inherently coincide with beat positions, making them natural targets for the pre-trained model’s rhythmic representations. Examining these results more closely, we verify that *Maracatu*’s tempo range of 165–180 BPM corresponds to inter-beat intervals of 333–363 ms. These durations approximately match the waveform spans of *cuica* and *gonge-lo*, but not those of the other instruments⁵.

This temporal alignment—where the instruments’ acoustic profile align with the genre’s inter-beat intervals—explains the high baseline accuracies. Additionally, the larger capacity of $TCNv2$ (116,302 parameters vs. 21,890 in $TCNv1$) and its exposure to a broader training set

⁵ According to an informal inspection of waveform spans—*cuica*: 384–428 ms, *gonge-lo*: 376–400 ms, *tarol*: 77–107 ms, *mineiro*: 90–180 ms, and *tambor-hi*: 120–230 ms.

may further contribute to this advantage. This suggests that model expressivity and pre-training diversity can compensate for task differences in certain transfer learning scenarios.

3.4 Discussion

Our investigation of two contrasting transfer learning scenarios reveals that adaptation outperforms baseline approaches across all instruments, with varying degrees of improvement.

In the within-domain setting, adaptation yielded high accuracies with F1 scores from 0.972 (*mineiro*) to 0.998 (*gongelo*) and 0.997 (*tarol*). Improvement was most pronounced for time-keeping instruments with lower baseline accuracies (≈ 0.500), with *cuica* showing a 52 p.p. gain. For the cross-domain adaptation, while improvements over the beat-tracking baseline ($_{\text{bsl}^*}$) were evident, comparison against the onset-tracking baseline ($_{\text{bsl}}$) revealed instrument-dependent patterns. Voicing instruments' best F1 scores remained below the onset-tracking baseline by 11-24 p.p., indicating limited benefits from domain adaptation. However, for time-keeping instruments whose onsets align with the pre-trained model's rhythmic priors, cross-task adaptation yielded improvements of 45-48 percentage points.

These findings provide key insights: i) Fine-tuning consistently enhances performance in both settings, making it valuable for achieving high accuracy in underrepresented music genres; ii) Models trained on beat-tracking can be effectively adapted for onset detection, leveraging model scale to compensate for task divergence and addressing limited data availability for non-mainstream instruments. However, effectiveness varies by instrument type: beat-to-onset adaptation benefits time-keeping instruments, while onset-to-onset adaptation consistently improves performance across all instruments. These improvements are naturally more substantial when baseline accuracy is lower, as observed in voicing instruments.

Our results also demonstrate that optimal fine-tuning configurations vary by instrument, necessitating tailored strategies for selecting which layer weights to update during fine-tuning. This challenges the assumption that only layers closest to the musical surface and the output layer would require recalibration to optimize a network for a specific instrument [17].

Finally, several limitations warrant consideration. Our results represent a single experimental cycle, and despite prior research suggesting relative stability across runs [28, 30], the stochastic nature of the (re)training process—due to convolutional dropout—implies that results may vary. While unlikely to affect general trends, multiple cycles would be needed to investigate specific aspects such as receptive field size impact and its relation to optimal layer freeze selection or instrument waveform profiles. Note that, as previously discussed, corresponding layers across the two models differ in their temporal receptive fields despite having the same labels. For instance, while $_{\text{Conv3}}$ corresponds to approximately 50 ms in both models, the

layer $_{\text{Tcn2}}$ spans 170 ms in $_{\text{TCNv1}}$ vs. 410 ms in $_{\text{TCNv2}}$. This discrepancy must be considered when interpreting results, limiting direct comparison between specific freeze configurations across scenarios. The lower annotation precision of *mineiro* further limits some result interpretation, potentially explaining its anomalous performance (e.g. lowest fine-tuned and baseline accuracy on each setting).

Notably, our current results were achieved with minimal adjustment to the experimental pipeline to maintain fair comparison with baselines. This conservative approach suggests greater improvements might be possible through hyperparameter optimization—for example, cross-task adaptation may require more epochs to converge than within-task adaptation. While such optimization exceeded this study's scope, it represents a promising direction for extending the clear performance gains demonstrated here.

4. CONCLUSION

This study investigated onset detection in *Maracatu de baque solto* through two transfer learning strategies: onset-to-onset adaptation and beat-to-onset adaptation. Both approaches yielded notable improvements over baseline models, underlining the advantages of fine-tuning for enhancing accuracy.

We demonstrated that cross-task adaptation of models is viable for less-represented tasks such as onset detection when structural alignment exists between source and target domains. Transfer learning effectively addresses limited data availability and circumvents extensive manual annotation or costly training from scratch—a finding with important implications for music information retrieval, particularly when facing data scarcity challenges.

Future work should address this study's limitations while exploring in greater detail the factors influencing transfer learning effectiveness. Multiple-run experiments would confirm observed trends and investigate specific aspects, such as optimal freeze segment selection and its relation with network receptive field and instrument waveform profiles, alongside potential improvements through hyperparameter optimization. Additional research directions include extending the analysis to other datasets and underrepresented instruments, and refining training protocols. Evaluating our adaptive approach using stricter tolerance windows would provide deeper insights into temporal precision, particularly for expressive instruments and microtiming analysis applications where fine-scale temporal variations are significant.

In summary, this study demonstrates the effectiveness of transfer learning in improving musical onset detection for diverse traditions beyond the Western canon. By adapting existing models, we can improve accuracy and robustness for underrepresented sounds. These methods and insights contribute to developing more inclusive tools for music analysis, with applications extending beyond the specific genres and tasks studied here to benefit the broader field of Music Information Retrieval.

5. REFERENCES

- [1] M. Goto and Y. Muraoka, "A beat tracking system for acoustic signals of music," in *Proceedings of the 2nd ACM International Conference on Multimedia (MULTIMEDIA '94)*. ACM Press, 1994, pp. 365–372.
- [2] S. Dixon, "Automatic Extraction of Tempo and Beat From Expressive Performances," *Journal of New Music Research*, vol. 30, no. 1, pp. 39–58, 2001.
- [3] R. B. Dannenberg, "Toward automated holistic beat tracking, music analysis, and understanding," in *Proceedings of the 6th International Conference on Music Information Retrieval (ISMIR)*, 2005, pp. 366–373.
- [4] S. Böck, F. Krebs, and M. Schedl, "Evaluating the on-line capabilities of onset detection methods," in *Proceedings of the 13th International Society for Music Information Retrieval Conference (ISMIR)*, 2012, pp. 49–54.
- [5] J. Pons, R. Gong, and X. Serra, "Score-informed syllable segmentation for a cappella singing voice with convolutional neural networks," in *Proceedings of the 18th International Society for Music Information Retrieval Conference (ISMIR)*, 2017, pp. 383–389.
- [6] R. Vogl, M. Dorfer, G. Widmer, and P. Knees, "Drum Transcription via Joint Beat and Drum Modeling using Convolutional Recurrent Neural Networks," in *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, 2017, pp. 150–157.
- [7] J. Bello, L. Daudet, S. Abdallah, C. Duxbury, M. Davies, and M. Sandler, "A tutorial on onset detection in music signals," *IEEE Transactions on Speech and Audio Processing*, vol. 13, no. 5, pp. 1035–1047, 2005.
- [8] S. Dixon, "Onset detection revisited," in *Proceedings of the 9th International Conference on Digital Audio Effects (DAFx)*, 2006, pp. 133–137.
- [9] M. Marolt, A. Kavcic, and M. Privosnik, "Neural networks for note onset detection in piano music," in *Proceedings of the International Computer Music Conference (ICMC)*, 2002.
- [10] A. Lacoste and D. Eck, "A Supervised Classification Algorithm for Note Onset Detection," *EURASIP Journal on Advances in Signal Processing*, vol. 2007, no. 1, p. 043745, 2006.
- [11] F. Eyben, S. Böck, B. Schuller, and A. Graves, "Universal onset detection with bidirectional long short-term memory neural networks," *Proceedings of the 11th International Society for Music Information Retrieval Conference, ISMIR 2010*, no. January, pp. 589–594, 2010.
- [12] J. Schlüter and S. Böck, "Musical onset detection with Convolutional Neural Networks," in *6th international workshop on machine learning and music (MML)*, 2013.
- [13] M. Tomczak and J. Hockman, "Onset Detection for String Instruments Using Bidirectional Temporal and Convolutional Recurrent Networks," in *Proceedings of the 18th International Audio Mostly Conference*. ACM, 2023, pp. 136–142.
- [14] G. Peeters, "The Deep Learning Revolution in MIR: The Pros and Cons, the Needs and the Challenges," in *Perception, Representations, Image, Sound, Music - 14th International Symposium, CMMR 2019, Marseille, France, October 14-18, 2019, Revised Selected Papers*, ser. Lecture Notes in Computer Science, R. Kronland-Martinet, S. Ystad, and M. Aramaki, Eds., vol. 12631. Springer, 2021, pp. 3–30.
- [15] A. Srinivasamurthy, A. Holzapfel, and X. Serra, "In Search of Automatic Rhythm Analysis Methods for Turkish and Indian Art Music," *Journal of New Music Research*, vol. 43, no. 1, pp. 94–114, 2014.
- [16] J. Bolt and G. Fazekas, "Supervised Contrastive Learning For Musical Onset Detection," in *Proceedings of the 18th International Audio Mostly Conference*. ACM, 2023, pp. 130–135.
- [17] M. E. P. Davies, M. Fuentes, J. Fonseca, L. Aly, M. Jerónimo, and F. B. Baraldi, "Moving in Time: Computational Analysis of Microtiming in Maracatu de Baque Solto," in *Proceedings of the 21th International Society for Music Information Retrieval Conference (ISMIR)*, 2020, pp. 795–802.
- [18] M. Fuentes, B. McFee, H. C. Crayencour, S. Essid, and J. P. Bello, "A Music Structure Informed Downbeat Tracking System Using Skip-chain Conditional Random Fields and Deep Learning," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, vol. 2019-May. IEEE, 2019, pp. 481–485.
- [19] A. Srinivasamurthy, A. Holzapfel, and X. Serra, "Informed automatic meter analysis of music recordings," in *Proceedings of the 18th International Society for Music Information Retrieval Conference (ISMIR)*, 2017, pp. 679–685.
- [20] D. Fioocchi, M. Buccoli, M. Zanoni, F. Antonacci, and A. Sarti, "Beat Tracking using Recurrent Neural Network: A Transfer Learning Approach," in *26th European Signal Processing Conference (EUSIPCO)*. IEEE, 2018, pp. 1915–1919.
- [21] G. Burloiu, "Interactive Learning of Microtiming in an Expressive Drum Machine," in *The Joint Conference on AI Music Creativity*, 2020.

- [22] Y. Wang, J. Salamon, M. Cartwright, N. J. Bryan, and J. P. Bello, “Few-Shot Drum Transcription in Polyphonic Music,” in *Proceedings of the 21st International Society for Music Information Retrieval Conference (ISMIR)*, 2020, pp. 117–124.
- [23] Y. Wang, N. J. Bryan, M. Cartwright, J. Pablo Bello, and J. Salamon, “Few-Shot Continual Learning for Audio Classification,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 321–325.
- [24] J. J. Valero-Mas and J. M. Iñesta, “Interactive user correction of automatically detected onsets: approach and evaluation,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2017, no. 1, p. 15, 2017.
- [25] K. Yamamoto, “Human-in-the-Loop Adaptation for Interactive Musical Beat Tracking,” in *Proceedings of the 22nd International Society for Music Information Retrieval Conference (ISMIR)*, 2021, pp. 794–801.
- [26] M. Schedl, E. Gómez, and J. Urbano, “Music Information Retrieval: Recent Developments and Applications,” *Foundations and Trends® in Information Retrieval*, vol. 8, no. 2-3, pp. 127–261, 2014.
- [27] A. S. Pinto and M. E. P. Davies, “Towards user-informed beat tracking of musical audio,” in *14th International Symposium on Computer Music Multidisciplinary Research (CMMR)*, 2019, pp. 577–588.
- [28] A. Pinto, S. Böck, J. Cardoso, and M. Davies, “User-Driven Fine-Tuning for Beat Tracking,” *Electronics*, vol. 10, no. 13, p. 1518, 2021.
- [29] L. S. Maia, M. Rocamora, and M. Fuentes, “Adapting meter tracking models to Latin american music,” in *Proceedings of the 23th International Society for Music Information Retrieval Conference (ISMIR)*, 2022, pp. 361–368.
- [30] A. S. Pinto and G. Bernardes, “Bridging the Rhythmic Gap : A User-Centric Approach to Beat Tracking in Challenging Music Signals,” in *16th International Symposium on Computer Music Multidisciplinary Research (CMMR)*, 2023, pp. 1–12.
- [31] K. Choi, G. Fazekas, M. Sandler, and K. Cho, “Transfer learning for music classification and regression tasks,” in *Proceedings of the 18th International Conference on Music Information Retrieval (ISMIR)*, 2017, pp. 141–149.
- [32] J. Fonseca, M. Fuentes, F. Bonini Baraldi, and M. E. Davies, “On the Use of Automatic Onset Detection for the Analysis of Maracatu de Baque Solto,” in *Perspectives on Music, Sound and Musicology. Current Research in Systematic Musicology*, vol. 10., L. Correia Castilho, R. Dias, and J. Pinho, Eds. Springer Cham, 2021, pp. 209–225.
- [33] S. Böck and M. E. P. Davies, “Deconstruct, Analyse, Reconstruct: How To Improve Tempo, Beat, and Downbeat Estimation,” in *Proceedings of the 21st International Society for Music Information Retrieval Conference (ISMIR)*, 2020, pp. 574–582.
- [34] B. McFee, J. W. Kim, M. Cartwright, J. Salamon, R. Bittner, J. P. Bello, and O.-s. Practices, “Open-Source Practices for Music Signal Processing Research: Recommendations for Transparent, Sustainable, and Reproducible Audio Research,” *IEEE Signal Processing Magazine*, vol. 36, no. January, pp. 128–137, 2019.
- [35] C. d. O. Santos, T. S. Resende, and P. M. Keays, *Batuque Book: Maracatu Baque Virado e Baque Solto*. Author’s edition., 2009.
- [36] G. P. Bessoni e Silva, “Maracatu de Baque Solto: de brincadeira a patrimônio cultural,” *Caderno Virtual de Turismo*, vol. 21, no. 2, p. 113, 2021.
- [37] M. E. P. Davies and S. Böck, “Temporal convolutional networks for musical audio beat tracking,” in *Proceedings of the 27th European Signal Processing Conference (EUSIPCO)*, 2019.
- [38] S. Böck, F. Korzeniowski, J. Schlüter, F. Krebs, and G. Widmer, “madmom: A New Python Audio and Music Signal Processing Library,” in *Proceedings of the 24th ACM International Conference on Multimedia (MM ’16)*. ACM, 2016, pp. 1174–1178.
- [39] J. Schlüter and S. Böck, “Improved musical onset detection with Convolutional Neural Networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 6979–6983.