

GAPS: A LARGE AND DIVERSE CLASSICAL GUITAR DATASET AND BENCHMARK TRANSCRIPTION MODEL

Xavier Riley Zixun Guo Drew Edwards Simon Dixon

Centre for Digital Music, Queen Mary University of London

j.x.riley@qmul.ac.uk, s.e.dixon@qmul.ac.uk

ABSTRACT

We introduce GAPS (Guitar-Aligned Performance Scores), a new dataset of classical guitar performances, and a benchmark guitar transcription model that achieves state-of-the-art performance on GuitarSet in both supervised and zero-shot settings. GAPS is the largest dataset of real guitar audio, containing 14 hours of freely available audio-score aligned pairs, recorded in diverse conditions by over 200 performers, together with high-resolution note-level MIDI alignments and performance videos. These enable us to train a state-of-the-art model for automatic transcription of solo guitar recordings which can generalise well to real world audio that is unseen during training.

For each track in the dataset, we provide metadata of the composer and performer, giving dates, nationality, gender and links to IMSLP or Wikipedia. We also analyse guitar-specific features of the dataset, such as the distribution of fret-string combinations and alternate tunings. This dataset has applications to various MIR tasks, including automatic music transcription, score following, performance analysis, generative music modelling and the study of expressive performance timing.

1. INTRODUCTION

Automatic Music Transcription (AMT) for instruments other than piano has faced challenges due to a lack of high-quality datasets [1]. This gap has limited the development of accurate transcription systems compared to those available for the piano, which benefit from comprehensive datasets like MAESTRO [2] and MAPS [3]. However, recent developments in audio-score alignment methods have shown promising results in improving transcription accuracy [1, 4].

With 2.7 million guitars sold in the US alone in 2019¹, the guitar is a popular instrument and retains a widespread cultural significance. Around 6% of these guitars sold were

of the classical or flamenco types (roughly 162,000 units). For comparison, around 31,000 acoustic pianos were sold in the US that year. Despite this popularity, we believe that the study of the guitar in the field of Music Information Retrieval (MIR) is underrepresented. Reviewing the paper titles for ISMIR conferences from 2013-2023 we find that publications with the word “piano” in the title outnumber those with “guitar” by 3 to 1². This imbalance may be due to the availability of high quality datasets for piano; new datasets and methods for guitar will help to address this.

In this paper, we present GAPS, a large and diverse classical guitar dataset that contains 14 hours of matched nylon string guitar audio recordings, note-level MIDI annotations, and corresponding music scores, where the recordings feature over 200 performers in diverse recording conditions. This is several times larger than GuitarSet [5], the EGDB dataset [6], the FrançoisLeduc dataset [4] and the IDMT-SMT-Guitar dataset [7] (see Section 2 for a detailed comparison). We use this data to train a benchmark transcription model which achieves state-of-the-art results for solo guitar transcription across 4 dataset splits.


The contributions of this paper are as follows:

- the largest available dataset consisting of real guitar audio, performance video, corresponding music scores and aligned MIDI annotations;
- metadata and external links for composers and performers, plus statistics of guitar-specific features;
- an efficient pipeline for verifying alignments of scores to audio;
- a benchmark state-of-the-art guitar transcription model trained on our dataset; and
- analysis and discussion of the effects of dataset quality, quantity and variety on AMT performance.

2. RELATED WORK

GuitarSet [5] is the most widely used MIR dataset for guitar. It provides around 3 hours of annotated guitar performances, where the data collection process required the use of a specialised guitar fitted with a hexaphonic pickup which was able to capture the output of individual strings. The use of a single guitar severely limits the diversity of

¹ <https://www.musictrades.com/us-retail-sales-guitar-market.html>

 © X. Riley, Z. Guo, D. Edwards and S. Dixon. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** X. Riley, Z. Guo, D. Edwards and S. Dixon, “GAPS: A Large and Diverse Classical Guitar Dataset and Benchmark Transcription Model”, in *Proc. of the 25th Int. Society for Music Information Retrieval Conf.*, San Francisco, United States, 2024.

² 46 piano and 15 guitar

Name	Audio type	Track count	Duration (m)	Note count	Scores
GuitarSet [5]	Real	360	180	62,476	No
IDMT-SMT-Guitar [7]	Real	1173	340	*5,767	No
EGDB [6]	Real	240	118	35,700	No
FrançoisLeduc [4]	Real	79	240	75,312	Yes (commercial)
GAPS (ours)	Real	300	843	259,410	Yes
SynthTab [8]	Synthetic	20,715	786,774	-	Yes, via DadaGP

Table 1. Comparison of existing guitar datasets, split into real and synthetic sources. * For IDMT, the note count is shown only for notes with annotations available.

timbres and recording conditions, and in turn makes it harder for AMT models to generalise from this data [4].

The EGDB [6] dataset contains 2 hours of guitar audio recorded by a professional guitarist using a hexaphonic pickup and recorded via DI (direct input). The DI signal is then further rendered using 6 different amplifier emulation plugins. The onsets and offsets of each note are annotated.

The IDMT-SMT-GUITAR database [7] is recorded by 3 musicians using 6 different guitars (5 electric, 1 acoustic). The final audio is either obtained from DI or microphone output. It contains 4 subsets each targeting a different MIR task, ranging from single notes to chords to various short musical pieces. Its utility in transcription tasks is limited however, as only a subset of the audio has corresponding time-aligned note annotations.

Improvements in diversity of audio sources were achieved by Maman and Bermano [1] through the use of score alignment techniques. Digital scores (in MIDI format) were aligned to the activations of the Onsets and Frames transcription model [9] trained on synthetic data. Low quality alignments were discarded and the remaining data was used to fine tune the model further. This expectation maximisation approach yielded a new state-of-the-art result on GuitarSet in the zero-shot setting, which demonstrated a generalisable model. The authors collected 5 hours of classical guitar recordings and scores in this work but these were not released as part of the publication.

Building on this approach, Riley et al. [4] published a new state-of-the-art model for guitar transcription. Instead of the Onsets and Frames model, they use the high resolution piano transcription model by Kong et al. [10], which was shown to be more tolerant of misaligned labels. Furthermore, instead of bootstrapping the process with synthetic data, they employ a pre-training step where a model is trained on the MAESTRO dataset with data augmentation, which was shown to improve generalisation. A dataset of around 4 hours of audio-MIDI pairs was published with their work, however the scores are not freely available as they were purchased from a commercial source.

As an alternative to annotating real world audio, Zang et al. [8] recently proposed a large scale dataset of synthesised audio from a subset of the DadaGP dataset [11]. When used as a pre-training step, the authors note improvements in multi-pitch estimation over 3 guitar datasets. De-

spite the large volume of additional training data, their note level results on the GuitarSet test split (86.1% F1 no off-set) are lower than those of several other methods which use GuitarSet alone (see [4]). This suggests that synthetic data alone is not sufficient to improve AMT systems, but a full comparison with consistent use of model architectures would be needed to establish this with certainty.

3. OVERVIEW OF DATASET

3.1 Dataset Curation

In an effort to improve the amount of available labelled, non-synthetic data, we have curated a new dataset of classical guitar recordings based on freely available scores from the ClassClef website³, together with matching performances on YouTube⁴. We align these sources using the automatic process described in [4] and then manually verified each alignment using the synchronised score viewer at soundslice.com. Following another alignment stage, any remaining scores with inaccurate alignments are rejected (using the criteria described below). This resulted in 300 performances sampled from the entire classical guitar canon totalling over 14 hours of music and over 250,000 note events. We have also curated extensive metadata, including information about the pieces, composers and performers, in order to enrich the dataset with details of the cultural context.

Our curation process is shown in Figure 1. It begins with the ClassClef website which provides around 5,500 pieces for download in PDF and GuitarPro formats. These focus mainly on the classical guitar with some flamenco and fingerstyle pieces included. Additionally, 547 of the pieces include links to videos on YouTube of a performance of the same piece. We first collected all GuitarPro files and converted them to MusicXML and MIDI formats using the free MuseScore software package⁵. We also downloaded the audio and video for the 547 pieces where YouTube links were available.

Using the alignment method described in [4], we produce an initial alignment between the score and the recording for each piece. This proceeds in two stages: an initial

³ classclef.com

⁴ youtube.com

⁵ musescore.org/en

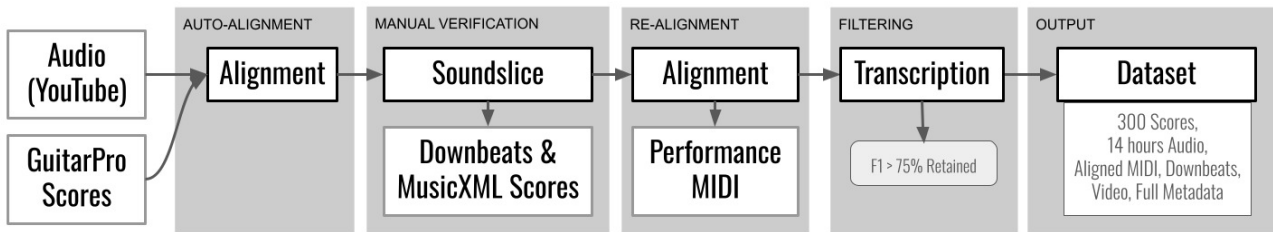


Figure 1. Flowchart of the dataset creation process.

alignment via Dynamic Time Warping (DTW), and a further fine alignment stage in which the notes of each chord are aligned to their closest activation from an existing transcription model. We emphasise this point as the resulting alignments are fully polyphonic in nature and as a result are more accurate than those produced by DTW alone, as described in [4].

In some cases the automatic alignment will not succeed, for example, when a linked video contains audio for an entire suite but the score only contains a single movement. For this reason a manual verification step was required. Using the `soundslice.com` website, we upload the automatically aligned downbeats to synchronize playback between the audio and the score. This allowed the authors of the paper (each with over 10 years of music experience) to review 474 of the scores (chosen at random) in an efficient workflow. More specifically, we manually verified the alignment between each downbeat location and the score for all 474 pieces. Particular attention was paid to the beginning and ending of each piece as these were a frequent source of issues in the DTW process. Moreover, any differences between the score and the performance that were identified were corrected, if feasible. In the end, 74 pieces were rejected for various reasons – for example those containing 7-string guitars, guitar duets and pieces where the edition did not match the performance. Out of the remaining 400 pieces examined, 280 were usable without corrections to the score and the remaining 120 required intervention to obtain correct downbeat alignments.

The 400 reviewed scores were then re-aligned using the same alignment method from step two of figure 1. The corrected downbeats were used as anchor points during this alignment stage to ensure that any alignment errors would be localised to one measure of music. To validate accuracy, we then compared our aligned versions of the score to outputs of the guitar transcription model from [4]. We retained the 300 scores with the highest agreement, measured using the “F-measure no offset” metric from the `mir_eval` library [12]. We retained scores which had an F-measure of more than 75%, yielding 300 audio-score pairs. We manually reviewed the lower scoring alignments and found a number of issues including errors with the processing of anacrusis bars, non-440Hz tunings and discrepancies between the performance and score editions. We hope to address these where possible as part of future work.

A summary of existing guitar datasets is shown in Table 1. When considering datasets with real (as opposed to syn-

thesised) audio, GAPS represents a significant advance in terms of the duration of audio and number of note events. In addition, ours is the first dataset of real audio to include freely available full music scores, tablatures in MusicXML format, and accompanying performance videos.

3.2 Composers

Works from 93 different composers are included, ranging from the Renaissance (Luys Milan, c.1500-1561) to the present day. The majority of works are from the classical guitar repertoire, with a small number of flamenco pieces and arrangements of popular music. We include the dates, nationality and presumed gender of each composer with links to canonical URLs (IMSLP and Wikipedia) where possible.

Examining the diversity of composers contained in the dataset, Figure 2 shows their nationalities, according to data from the canonical URL for each composer. This shows they are broadly divided between Europe and Latin America. In terms of chronology Figure 3 shows the distribution of pieces according to the year in which the composer was born. This shows that the included pieces are mainly weighted around the Romantic era (1850-1900). The peak around 1650 is almost entirely due to J.S. Bach, who is the second most common composer in our dataset with 23 pieces. We also include information about the presumed gender of composers in our metadata, however only two female composers (Maria Linnemann and Luise Walker) are included who together represent 2% of the total by piece count. We acknowledge that this is a shortcoming of the current dataset and we will seek to address this in future work.

3.3 Performances

The accompanying videos are drawn from 205 different performers with YouTube views totalling over 35 million across all videos. Some are professionally produced recordings whereas others are recorded on commodity equipment such as phones and laptops. We believe this is an advantage of this dataset in that recordings are drawn from a wide variety of real world recording conditions, which in turn helps to increase the robustness of trained AMT models.

In the metadata we include information about the name of the performer (where available), their social media links (if available), the YouTube channel, the view count and

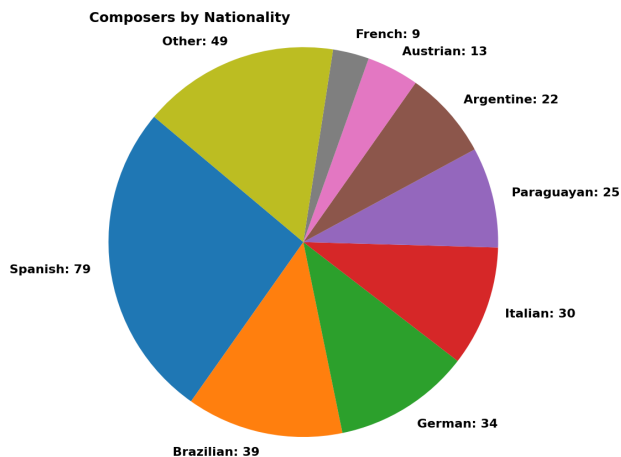


Figure 2. Nationalities of the composers

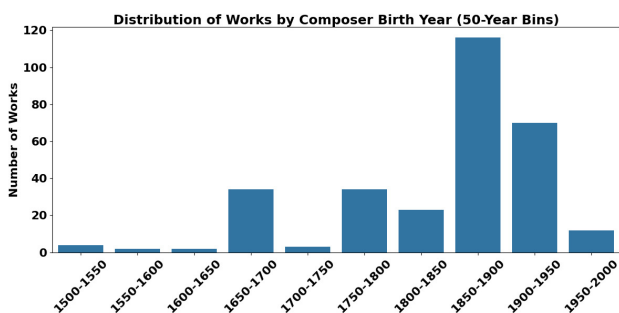


Figure 3. Histogram of works according to composer’s birth year at 50 year intervals

Tuning	Count	% of total
EADGBE	232	77.33
DADGBE	58	19.33
DGDGBE	5	1.67
EADF#BE	2	0.67
FADGBE	1	0.33
CGDGBE	1	0.33
EBDGBE	1	0.33

Table 2. Distribution of guitar tunings in GAPS. The tuning is expressed from low to high pitch.

the presumed gender of the performer. This was gathered to examine the extent to which classical guitar is a male dominated field. We find that female performers are better represented than composers in our dataset, but still only comprise 23% of the total.

3.4 Guitar-Specific Features

The large number of scores allows us to examine several guitar-specific features of the data. In Table 2 we see that two different tunings account for 97% of the data. While standard tuning is most common, almost 20% of pieces have the lowest string tuned down one tone to D. Other alternate tunings account for around 3.3% of the total.

To see the distribution of notes across the guitar neck in

this dataset, we have plotted a heat map as shown in Figure 4 using the fret information contained in the MusicXML tablature. Over the 259,000 note events we see that the pieces in the classical guitar repertoire favour the use of open strings and the first position. The strong peak at the 2nd fret A on the G string also suggests a preference towards “guitar friendly” keys such as E and A which allow the performer to use the open bass and top strings. While this distribution is uneven, we consider this to be representative of the classical guitar repertoire. We encourage other dataset authors to explore similar visualisations in future work to see if this varies with other genres.

Since most pitches can be played on more than one position on the guitar, there is an exponentially large number of tablatures that correspond to any one given score, including many physically unplayable versions. While each tablature in our dataset represents one valid way to play the score, we have not verified the extent to which the tablatures correspond to the choices of the performers in the specific performances in the GAPS dataset. This is left for future work. As we were not able to trace the provenance of the ClassClef data, we presume the data is crowdsourced and reflects the playing habits of a subset of computer-literate guitarists. It is also possible that some of the tabs were generated algorithmically from the score data.

4. TRANSCRIPTION BASELINE

4.1 Experimental Settings

To demonstrate the utility of the GAPS dataset of aligned score-audio pairs, we trained several guitar transcription models using the high resolution model of Kong et al. [10], which achieved state-of-the-art performance when trained for guitar transcription [4]. This model is a convolutional recurrent neural network (CRNN) that is trained in a supervised manner to map log mel-spectrograms of 10-second segments of audio to MIDI. The convolutional layers span only across the frequency dimension, maintaining the time-resolution of the original spectrogram (10ms). These features are then processed by a gated recurrent unit (GRU) to produce the final outputs of onset, offset, frame activity, and velocity activations per pitch per time window.

There are two reasons why we used the high resolution model [10]. Firstly to ensure fair comparisons with the state-of-the-art model in [4] as it shares the same architecture. This allows us to examine how our GAPS dataset influences the same transcription model. Secondly, fine-tuning becomes feasible due to the shared architecture among multiple piano transcription models [10, 13]. This allows us to investigate whether different pre-trained piano transcription models can improve guitar transcription through domain adaptation.

For our experiments, we trained 2 sets of models. The first set of models is trained only on the GAPS dataset and the second set of models is trained with a combination of GuitarSet and GAPS. We employ the first set of models for zero-shot inference on the complete GuitarSet, while the

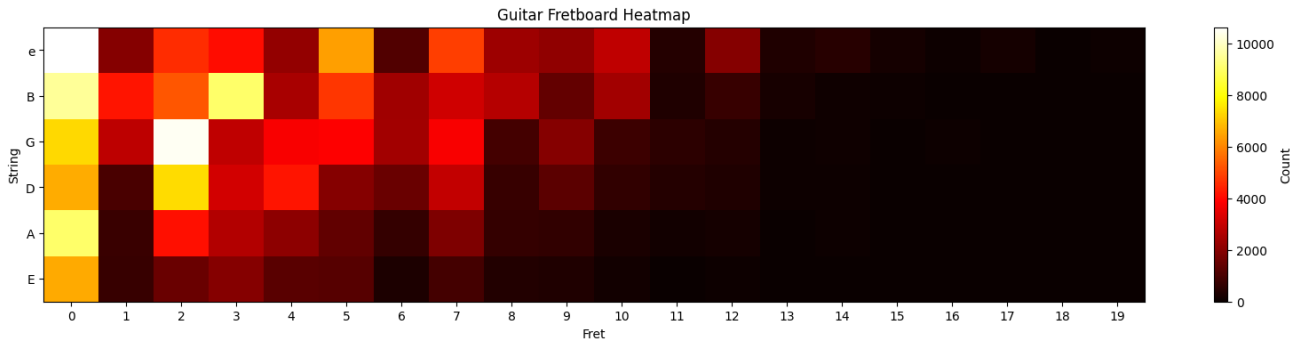


Figure 4. Heat map of the fret/string combinations in the GAPS MusicXML tablatures.

second set is utilised to evaluate guitar transcription performance across the test splits of GuitarSet, theFrançoisLeduc dataset and GAPS. To study the effects of pre-training and finetuning [8, 13], each set of models has 3 variants: one trained from scratch and two finetuned from one of two published checkpoints for piano transcription [10, 13]. This also allows for a more direct comparison with results reported in [4].

Regarding our training data and strategy, we randomly divide the GAPS dataset with a 90:10 split by piece, for training and testing respectively. Following [4], each audio file is split into 10-second chunks, using a hop size of 1 second. We adopt the same train-test split from [4, 14] for GuitarSet. During training, pitch shifting of up to ± 3 semitones was randomly applied as data augmentation [14].

4.2 Transcription Results

In Tables 3 to 6, we report the evaluation results for the models described in Section 4.1. Our proposed combination of model, pre-training checkpoint and dataset achieves state-of-the-art performances on all 4 test sets mentioned in Section 4.1. Considering the similarities to the approach used by Riley et al. [4], our larger dataset appears to drive the improvement in results.

4.2.1 Generalisation and Guitar Types

GuitarSet contains audio for one acoustic steel string guitar recorded via microphone and also via the guitar pickup (the “DI” outputs). Despite our GAPS data containing only performances on nylon-stringed classical guitars, our model is able to generalise well to GuitarSet in the zero-shot setting (F-measure 88.1% - see Table 4). This result is interesting as it appears that timbral differences between guitars are not a strong factor in the success of the model for this task. On the other hand, GAPS does include a large range of guitars and recording conditions (unlike GuitarSet’s one guitar), which we expect would contribute to the generalisation performance of models trained on it.

We also note that for the other solutions based on encoder-decoder architectures [14, 16], the strong results in the supervised setting on GuitarSet fail to perform as well on unseen data. Table 4 shows the transcription accuracy on GuitarSet in the zero-shot setting, i.e. where models are trained without any access to GuitarSet. F-measure scores

for MT3 fall from 90.0% to 32.0% on GuitarSet. The previous state-of-the-art model (Time-Frequency Perceiver) from Lu et al. [14] attains 91.1% in the GuitarSet supervised task but drops to 80.0% on the unseen FrançoisLeduc test set. It may be the case that these architectures require more data to generalise effectively and we hope to explore training them on GAPS in future work.

For the FrançoisLeduc test split in Table 5, our proposed model outperforms Riley et al. [4] by a small margin, however their model was trained in a supervised fashion whereas this dataset was unseen by our model.

Conversely, our proposed method outperforms Riley et al. [4] on the GAPS test split by a margin of 2.2% (see Table 6). This indicates that, despite our method’s strong generalisation (see Table 4), it is somewhat specialised to classical guitar timbres and that the strongest results in the future may rely on the use of specific training data.

	P_{50}	R_{50}	F_{50}
Basic Pitch [15]	-	-	79.0
MT3 [16]	-	-	90.0
Zang et al [8]	-	-	84.5
Lu et al. [14]	-	-	91.1
SpecTNT (in [14])	-	-	90.7
Riley et al. [4] (FL)	87.6	86.8	86.9
Riley et al. ($GS+FL$)	91.1	88.5	89.7
Ours			
($GAPS$)	89.9	85.4	87.2
($GAPS$ Finetuned from [10])	88.8	86.8	87.5
($GAPS$ Finetuned from [13])	90.1	86.6	88.0
($GAPS+GS$)	90.2	90.9	90.4
($GAPS+GS$ Finetuned from [10])	89.4	92.1	90.7
($GAPS+GS$) Finetuned from [13])	91.3	90.7	91.2

Table 3. Results for note-level transcription accuracy on the GuitarSet test split. P_{50} , R_{50} , and F_{50} are Precision, Recall and F1-measure, expressed as percentages, at 50ms resolution. All are evaluated on onsets only (no offsets or velocity), using the `mir_eval` library. Baseline results are described in [4].

	P_{50}	R_{50}	F_{50}
MT3 [16]	-	-	32.0
Kong et al. [10]	67.5	49.7	54.8
Kong et al. (w/ aug)	80.6	44.0	50.3
Zang et al. [8] (Synthtab)	-	-	70.2
Maman (MusicNet _{EM}) [1]	86.6	80.4	82.9
Maman (Guitar) [1]	86.7	79.7	82.2
Riley et al. [4]	88.0	87.1	87.3
Ours	92.4	81.8	86.1
Ours (Finetuned from [10])	91.6	83.7	87.0
Ours (Finetuned from [13])	91.1	85.9	88.1

Table 4. Results for note-level transcription accuracy on the entire GuitarSet in the zero-shot setting.

	P_{50}	R_{50}	F_{50}
Basic Pitch [15]	54.6	85.0	66.1
Omnizart [17]	63.0	72.1	67.1
MT3 [16]	48.8	57.0	52.4
Lu et al. [14]	83.6	77.3	80.0
Riley et al. [4]	83.9	85.5	84.7
Ours (Finetuned from [13])	85.5	84.2	84.8

Table 5. Results for note-level transcription accuracy on the test split of the FrançoisLeduc dataset [4].

	P_{50}	R_{50}	F_{50}
Riley et al. [4]	92.9	91.4	92.1
Ours	94.9	92.1	93.4
Ours (Finetuned from [10])	94.6	93.4	94.0
Ours (Finetuned from [13])	95.0	93.6	94.3

Table 6. Results for note-level transcription accuracy on a test split of the GAPS dataset.

4.2.2 Effects of Pre-training

In each of our evaluations, we see a consistent trend whereby the model with no pre-training is surpassed by the model pre-trained on piano (MAESTRO) and fine-tuned on GAPS, which in turn is surpassed by the model pre-trained on an augmented version of MAESTRO [13] before fine-tuning on GAPS. This illustrates the importance of pre-training and fine-tuning, as well as data augmentation as important drivers of success in the transcription task (see Edwards et al. [13] for a detailed analysis of the effect of data augmentation on transcription generalisation).

We also note that strong results for other methods on the GuitarSet test split are obtained from models trained with a mixture of datasets [4, 14, 16]. One exception is Zang et al. [8], who use a large corpus of synthetically rendered guitar samples for pre-training. This does not perform as well as other methods but their results were obtained from a model (TabCNN) designed for guitar tablature prediction, as opposed to a state-of-the-art transcription model. A full comparison of synthetic and real audio for pre-training is something we also hope to explore in future work.

5. CONCLUSION

We present GAPS, a large dataset of score and audio pairs for solo classical guitar which comprises a wide range of composers, performers and real-world recording conditions, totaling 14 hours of recordings. The MIDI annotations are made freely available and the audio is available at the YouTube links provided. This represents the largest dataset of freely available guitar audio-score pairs to date.

We included analysis of the overall statistics of the GAPS dataset, but further musicological work could be done to examine connections between the composers, performers and musical features. The published MIDI annotations could be useful for generative modelling of classical guitar and other instruments. For future work we will look to expand the dataset and enhance the diversity where possible, particularly for the range of composers we include.

One application of this dataset is AMT for guitar, which we demonstrate through a comprehensive evaluation of a transcription model trained on our data. This shows state-of-the-art results when compared with existing methods trained on other datasets. In future work we look to examine further issues around pre-training for guitar transcription.

6. ETHICS STATEMENT

In addition to our role as researchers, we are also members of the global community of musicians and we seek to respect their important role in our culture. Our work here raises several issues which may have wider impact on this community which we hope to address as follows.

Firstly, we believe that using sources which are publicly available (subject to licence conditions) is important to reduce barriers to future research. At the time of writing, neither the scores nor their audio recordings are behind any kind of paywall. We have processed this data and make the results available on the basis of fostering research. We also obtained permission from the website owner of `classclef.com` to make use of their materials.

By publishing work on YouTube, artists do grant some kind of implicit licence that the data can be viewed, however the specific terms of the licence may restrict further use cases. We believe that our work is justified in using this data under fair use or fair dealing exemptions defined for research, but we are mindful that further use of the data may require express permission from the performers, composers or copyright-holders. We have attempted to address this by including detailed information about all performers and composers in the accompanying metadata to allow interested parties to contact them directly.

Finally we recognise that AMT models which approach human-level accuracy might pose a threat to those who are employed in music transcription and related fields. On the other hand, such models could also assist such work and become tools for improving the efficiency and accuracy of their daily work. For this reason we are carefully considering whether to make our model weights freely available.

7. ACKNOWLEDGMENTS

Authors XR, ZG and DE are research students at the UKRI Centre for Doctoral Training in Artificial Intelligence and Music, supported by UK Research and Innovation [grant number EP/S022694/1] and Yamaha Corporation (DE).

8. REFERENCES

- [1] B. Maman and A. H. Bermann, “Unaligned supervision for automatic music transcription in the wild,” in *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, ser. Proceedings of Machine Learning Research, vol. 162. PMLR, 2022, pp. 14 918–14 934.
- [2] C. Hawthorne, A. Stasyuk, A. Roberts, I. Simon, C. A. Huang, S. Dieleman, E. Elsen, J. H. Engel, and D. Eck, “Enabling factorized piano music modeling and generation with the MAESTRO dataset,” in *7th International Conference on Learning Representations*, New Orleans, USA, 2019.
- [3] V. Emiya, N. Bertin, B. David, and R. Badeau, “MAPS - a piano database for multipitch estimation and automatic transcription of music,” INRIA, France, Research Report 00544155, 2010. [Online]. Available: <https://hal.inria.fr/inria-00544155>
- [4] X. Riley, D. Edwards, and S. Dixon, “High resolution guitar transcription via domain adaptation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, 2024*. IEEE, 2024. [Online]. Available: <https://arxiv.org/abs/2402.15258>
- [5] Q. Xi, R. M. Bittner, J. Pauwels, X. Ye, and J. P. Bello, “GuitarSet: A dataset for guitar transcription,” in *Proceedings of the 19th International Society for Music Information Retrieval Conference*, Paris, France, 2018, pp. 453–460.
- [6] Y. Chen, W. Hsiao, T. Hsieh, J. R. Jang, and Y. Yang, “Towards automatic transcription of polyphonic electric guitar music: A new dataset and a multi-loss transformer model,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022*. IEEE, 2022, pp. 786–790.
- [7] C. Kehling, J. Abeßer, C. Dittmar, and G. Schuller, “Automatic tablature transcription of electric guitar recordings by estimation of score- and instrument-related parameters,” in *Proceedings of the 17th International Conference on Digital Audio Effects, DAFX-14, Erlangen, Germany, September 1-5, 2014*, 2014, pp. 219–226.
- [8] Y. Zang, Y. Zhong, F. Cwitkowitz, and Z. Duan, “Synthtab: Leveraging synthesized data for guitar tablature transcription,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, 2024*. IEEE, 2024. [Online]. Available: <https://arxiv.org/abs/2402.15258>
- [9] C. Hawthorne, E. Elsen, J. Song, A. Roberts, I. Simon, C. Raffel, J. H. Engel, S. Oore, and D. Eck, “Onsets and frames: Dual-objective piano transcription,” in *Proceedings of the 19th International Society for Music Information Retrieval Conference*, Paris, France, 2018, pp. 50–57.
- [10] Q. Kong, B. Li, X. Song, Y. Wan, and Y. Wang, “High-resolution piano transcription with pedals by regressing onset and offset times,” *IEEE ACM Transactions on Audio, Speech and Language Processing*, vol. 29, pp. 3707–3717, 2021.
- [11] P. Sarmento, A. Kumar, C. J. Carr, Z. Zukowski, M. Barthelet, and Y. Yang, “Dadagp: A dataset of tokenized guitarpro songs for sequence models,” in *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, 2021, pp. 610–617.
- [12] C. Raffel, B. McFee, E. J. Humphrey, J. Salamon, O. Nieto, D. Liang, and D. P. W. Ellis, “Mir_eval: A transparent implementation of common MIR metrics,” in *Proceedings of the 15th International Society for Music Information Retrieval Conference*, Taipei, Taiwan, 2014, pp. 367–372.
- [13] D. Edwards, S. Dixon, E. Benetos, A. Maezawa, and Y. Kusaka, “A data-driven analysis of robust automatic piano transcription,” *IEEE Signal Process. Lett.*, vol. 31, pp. 681–685, 2024.
- [14] W. T. Lu, J. Wang, and Y. Hung, “Multitrack music transcription with a time-frequency perceiver,” in *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*. IEEE, 2023, pp. 1–5.
- [15] R. M. Bittner, J. J. Bosch, D. Rubinstein, G. Meseguer-Brocal, and S. Ewert, “A lightweight instrument-agnostic model for polyphonic note transcription and multipitch estimation,” in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Singapore, 2022.
- [16] J. Gardner, I. Simon, E. Manilow, C. Hawthorne, and J. H. Engel, “MT3: Multi-task multitrack music transcription,” in *Tenth International Conference on Learning Representations*, 2022.
- [17] Y. Wu, Y. Luo, T. Chen, I. Wei, J. Hsu, Y. Chuang, and L. Su, “Omnizart: A general toolbox for automatic music transcription,” *J. Open Source Softw.*, vol. 6, no. 68, p. 3391, 2021. [Online]. Available: <https://doi.org/10.21105/joss.03391>