



ReconVAT: A Semi-Supervised Automatic Music Transcription Framework for Low-Resource Real-World Data

Kin Wai Cheuk*

Dorien Herremans

Information Systems Technology and Design
Singapore University of Technology and Design
Singapore

{kinwai_cheuk,dorien_herremans}@mymail.sutd.edu.sg

ABSTRACT

Most of the current supervised automatic music transcription (AMT) models lack the ability to generalize. This means that they have trouble transcribing real-world music recordings from diverse musical genres that are not presented in the labelled training data. In this paper, we propose a semi-supervised framework, ReconVAT, which solves this issue by leveraging the huge amount of available unlabelled music recordings. The proposed ReconVAT uses reconstruction loss and virtual adversarial training. When combined with existing U-net models for AMT, ReconVAT achieves competitive results on common benchmark datasets such as MAPS and MusicNet. For example, in the few-shot setting for the string part version of MusicNet, ReconVAT achieves F1-scores of 61.0% and 41.6% for the note-wise and note-with-offset-wise metrics respectively, which translates into an improvement of 22.2% and 62.5% compared to the supervised baseline model. Our proposed framework also demonstrates the potential of continual learning on new data, which could be useful in real-world applications whereby new data is constantly available.

CCS CONCEPTS

- Applied computing → Sound and music computing;
- Computing methodologies → Semi-supervised learning settings; Neural networks.

KEYWORDS

semi-supervised training, virtual adversarial training, audio processing, automatic music transcription, music information retrieval

ACM Reference Format:

Kin Wai Cheuk, Dorien Herremans, and Li Su. 2021. ReconVAT: A Semi-Supervised Automatic Music Transcription Framework for Low-Resource Real-World Data. In *Proceedings of the 29th ACM International Conference on Multimedia (MM '21), October 20–24, 2021, Virtual Event, China*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3474085.3475405>

*Also with Institute of High Performance Computing (IHPC), Agency for Science, Technology and Research (A*STAR).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '21, October 20–24, 2021, Virtual Event, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3475405>

Li Su

Institute of Information Science
Academia Sinica
Taiwan
lisu@iis.sinica.edu.tw

1 INTRODUCTION

Automatic Music Transcription (AMT), is a fundamental problem in the field of Music Information Retrieval (MIR). According to the definition from the field of Music Information Retrieval (MIR) [3], AMT aims at transcribing music audio files into symbolic representations such as piano rolls [7] or music scores [6, 38, 39], which is very similar to Automatic Speech Recognition (ASR) [1, 8]. These symbolic representations have a wide range of applications including music indexing [15, 42], music generation [20], music recommendation system (MRS) [9, 32, 48], music analysis [22, 26, 30], and automatic music accompaniment [31].

Recent advances in fully supervised deep learning have enabled AMT models [18, 19, 23, 24] to achieve state-of-the-art performance for solo piano pieces, given sufficient labelled training data. While acoustic audio recordings as well as the aligned midi labels for piano music can be easily obtained by using a hybrid acoustic/midi piano such as the Yamaha Disklavier [16], this is not the case for other musical instruments such as violin and clarinet. At the time of writing, hybrid versions of these musical instruments are still not available. They are either midi controllers that lack the capability to produce realistic acoustic sound, or fully acoustic instruments without the capability to record the midi annotation in real-time. Therefore, the paired acoustic recordings and midi annotations for these instruments are very expensive to obtain, and hence, very limited. Supervised models fail to function well for these instruments.

Self-supervised or semi-supervised learning is an underexplored area in AMT. Existing unsupervised models have only been applied to specific musical instruments. For example, Berg-Kirkpatrick et al. [5] proposed an unsupervised graphical model using prerecorded key-wise piano samples to reconstruct the original signal. Upon successful reconstruction, the model could infer the transcription result via both the predicted onset locations and the selected piano samples for reconstructing the spectrogram. Choi and Cho [14] also employed a similar approach in their unsupervised drum transcription model. This approach, however, only works when the musical instrument is a percussion or plucked instrument type with a clear transient, immediately followed by a natural decay (piano is considered as a percussion instrument due to the hammering mechanism). Musical instruments which produce an increasing or fluctuating amplitude after the transient are unable to be represented by a fixed audio sample, and hence cannot be properly transcribed using the above-mentioned approach. Examples of such instruments include string instruments that are capable of starting a long note softly

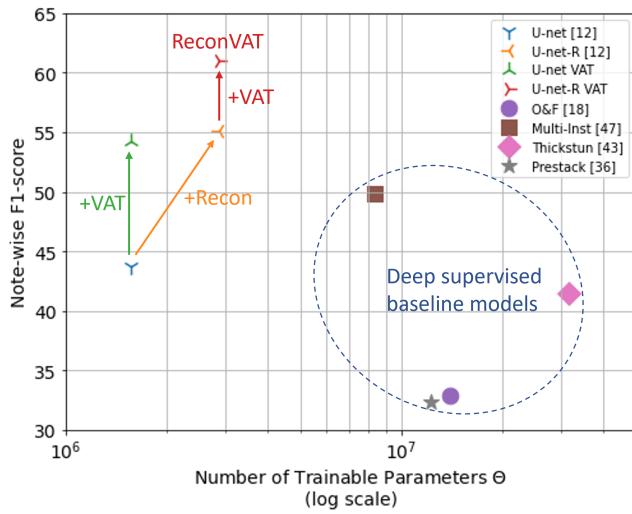


Figure 1: A scatter plot showing the note-wise F1-score and the number of model parameters for different models trained on the string version of MusicNet (see Section 4.2 for details).

followed by a crescendo through gradually increasing the bow pressure; and woodwind instruments that can sustain a note as long as the player’s lung capacity can handle.

In this paper, we propose a semi-supervised AMT framework, which we refer to as **ReconVAT**, based on the idea of spectrogram reconstruction [12] and virtual adversarial training (VAT) [34]. ReconVAT leverages unlabelled data to improve the transcription accuracy with only a limited amount of labelled data, and it works well on various musical instruments such as piano, string instruments, as well as woodwind instruments. We also show that our framework has important applications such as continual learning with new unlabelled recordings, and being able to transcribe music genres that are outside of the labelled training set. More importantly, all of these can be achieved with only a small number of model parameters compared to existing deep learning models for AMT as shown in Figure 1. This makes our framework attractive and practically usable for real-world applications deployed on mobile devices. To the best of our knowledge, this is the first semi-supervised deep learning framework for instrument-agnostic AMT at the time of writing.

The contributions of this paper can be summarized as follows:

- We propose a semi-supervised framework for AMT that generalizes well across *different kinds of musical instruments*.
- We leverage existing models by integrating them into the proposed semi-supervised framework to achieve state-of-the-art transcription accuracy for low-resource scenario.
- We demonstrate possible applications in continual learning on music genres that are not present in the train set.

2 METHOD

In this section, we will formulate automatic music transcription (AMT) mathematically. Then we will introduce related work, and

describe how to combine both spectrogram reconstruction [12] and VAT [34] to be our proposed semi-supervised framework ReconVAT.

2.1 Problem Definition

The goal of AMT is to convert audio data into symbolic music data [3, 4]. In this paper, we consider the case of converting spectrograms into piano rolls. Given a normalized input spectrogram $X_{\text{spec}} \in [0, 1]^{T \times F}$, where T is the number of timesteps and F is the number of frequency bins, we want to have a model $p(Y_{\text{post}}|X_{\text{spec}}, \theta)$, with a set of trainable parameters θ , that infers the posteriorgram $Y_{\text{post}} \in [0, 1]^{T \times N}$. Here N is the note range for the musical instrument, for example, $N = 88$ for piano transcription since there are 88 keys on the keyboard. The ground truth piano roll $Y_{\text{roll}} \in \{0, 1\}^{T \times N}$ is the symbolic notation we want to predict. This is done by simply applying a threshold (e.g. 0.5) to the Y_{post} .

2.2 Spectrogram Reconstruction

Cheuk et al. [12] proposed a model consisting of a transcriber $p(Y_{\text{post}}|X_{\text{spec}}, \theta)$ and a reconstructor $q(X_{\text{recon}}|Y_{\text{post}}, \phi)$. The reconstructor uses the posteriorgram generated from the transcriber as input to reconstruct the spectrograms X_{recon} . Therefore, in addition to the transcription loss $L_{\text{trans}}(Y_{\text{post}}, Y_{\text{label}})$, there is also a reconstruction loss $L_{\text{recon}}(X_{\text{recon}}, X_{\text{spec}})$ to be minimized.

The reconstructed spectrograms X_{recon} are then used to train the same transcriber $p(Y'_{\text{post}}|X_{\text{recon}}, \theta)$ again, resulting in one extra transcription loss $L_{\text{trans}}(Y'_{\text{post}}, Y_{\text{label}})$. Cheuk et al. [12] has shown that training the model in this manner results in a consistently better model. Their reported results, however, are unable to beat the state-of-the-art AMT models. We will show in Table 1 that their model can be modified to compete with the state-of-the-art AMT models. Although it is not demonstrated in their paper, they also claim that their model has the potential to be trained in an unsupervised manner. We will therefore also show in Section 2.4 that when combined with virtual adversarial training [34], we can modify the spectrogram reconstruction framework [12] to be a semi-supervised model for AMT.

2.3 Virtual Adversarial Training

Virtual adversarial training (VAT), as presented by Miyato et al. [34] is an extended version of adversarial training (AT) proposed by Goodfellow et al. [17]. In AT, labels are required to calculate the adversarial vectors. In the case where we do not have access to the labels, Miyato et al. [34] proved that the adversarial vector can be obtained via Equation (1):

$$r_{\text{adv}}^{\text{VAT}} = \epsilon \left(\nabla_r D \left[p(Y_{\text{pred}}|X, \theta), p(Y_{\text{adv}}|X + r, \hat{\theta}) \right] \right), \quad (1)$$

where r is a randomly initialized noise vector, and Y_{adv} is the output obtained using the adversarial input $X_{\text{adv}} = X + r$.

By doing so, it is possible to perform adversarial training using unlabelled datasets. Most existing literature applies VAT to static classifications [1, 8, 27, 28, 33, 35]. While SeqVAT [10] is designed for sequential labelling, it is a one-hot prediction system. To the best of our knowledge, ReconVAT is the first framework capable of multi-hot sequential labelling for polyphonic AMT. In the next section,

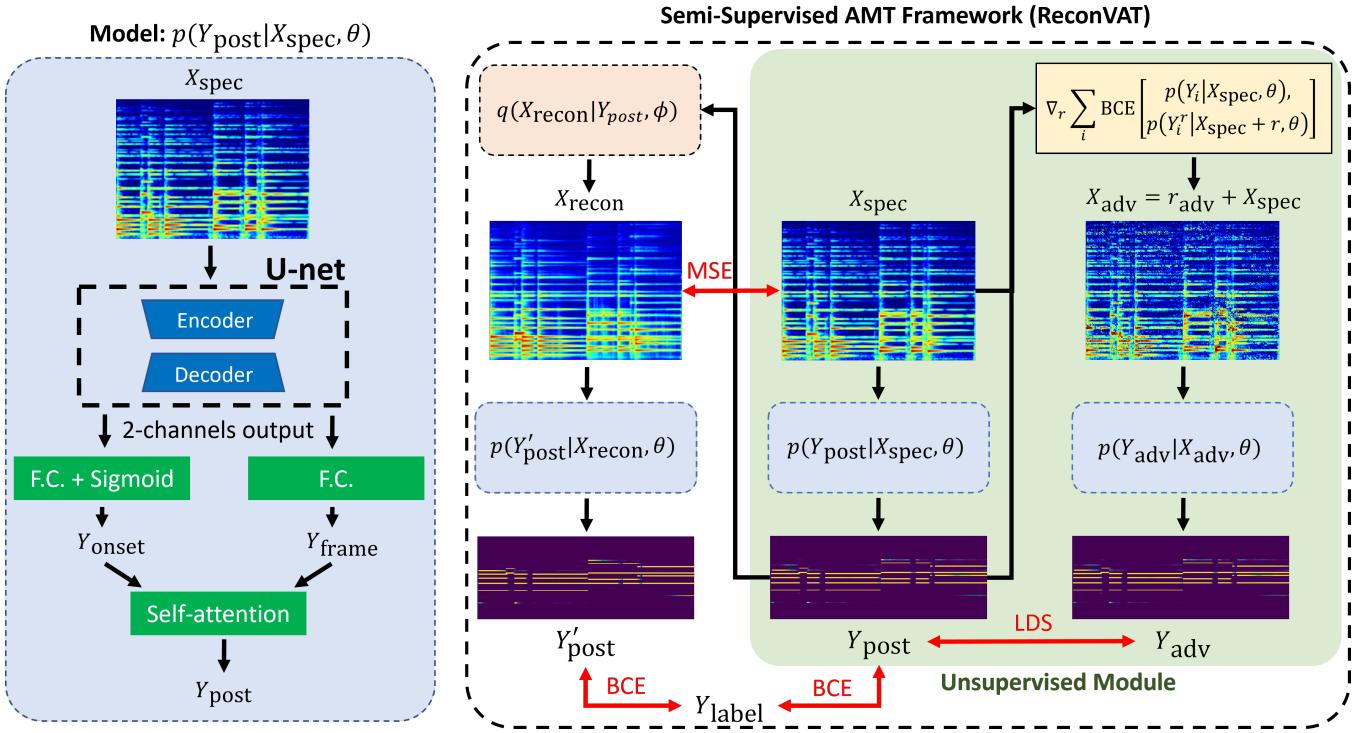


Figure 2: The left-hand side of the figure shows the modified version of the AMT model proposed in [12] such that it supports onset prediction. The right-hand side of the figure shows the proposed framework for semi-supervised AMT. For simplicity, we omit the onset prediction in the figure by showing the case when $i = \{\text{post}\}$. The region highlighted in green is the unsupervised module which supports training not only on labelled samples, but also on unlabelled samples.

we will describe how to combine both spectrogram reconstruction and VAT to obtain a semi-supervised framework for AMT.

2.4 Proposed Framework – ReconVAT

The left-hand side of Figure 2 shows our modified version of the model proposed by Cheuk et al. [12]. We improve it by introducing a two-channel output, one channel for the onset prediction Y_{onset} , and another channel for the frame feature extraction Y_{frame} . The posteriorogram is obtained from a self-attention layer which takes the concatenation of Y_{onset} and Y_{frame} as the input. This modification increases the model’s flexibility. For example, if we want to also include the offset prediction Y_{offset} , we can have a three-channel output instead. For simplicity, we will only explore the case of one-channel (Y_{frame}) and two-channel (Y_{frame} and Y_{onset}) prediction in this paper. The implementation details will be discussed in Section 3.4.

The right-hand side of Figure 2 shows our proposed ReconVAT. It consists of three branches. The framework starts with the middle branch where it takes X_{spec} as the input and outputs a posteriorogram Y_{post} . The branch on the left then takes the Y_{post} as its input and generates a reconstructed spectrogram X_{recon} using the reconstructor $q(X_{\text{recon}}|Y_{\text{post}}, \phi)$ mentioned in 2.2. The reconstructed spectrogram X_{recon} is passed to the same model again to obtain another posteriorogram Y'_{post} . The two posteriorograms Y_{post} and Y'_{post} should be as close to the label Y_{label} as possible.

The branch on the right-hand side is the unsupervised module which uses VAT. To obtain the adversarial spectrogram X_{adv} , we apply a modified version of VAT that works better for AMT (Section 3.4). Using this adversarial spectrogram, we obtain another posteriorogram Y_{adv} via the same model.

For labelled spectrograms, all three branches are used. For unlabelled spectrograms, we only use the middle and the right branches (highlighted in green in Figure 2). This framework is trained by minimizing both the supervised loss L_1 (Equation 5) and the unsupervised loss L_{ul} which will be discussed in detail in Section 3.5.

3 EXPERIMENTS

In this section, we describe the datasets and the experiments for demonstrating the power of our proposed semi-supervised framework ReconVAT.

3.1 MAPS dataset

The MAPS dataset [16] consists of nine folders, each folder contains 30 full-length midi recordings. In seven of these folders, the audio recordings are synthesized from the midi annotations using different virtual piano software such as Steinberg, Native Instruments, and Sampletek. Only in the folders ENSTDkAm and ENSTDkC1, the audio recordings are recorded simultaneously with the midi recordings using a Yamaha Disklavier. We follow the existing consensus [12, 18, 23, 36, 41] that the seven folders containing artificially

generated audio recordings should be used as the training set, and the other two folders, ENSTDkAm and ENSTDkC1, as the test set.

Since some music pieces appear in both the training and the test set, we follow the existing literature [18, 41] to remove overlapping songs from the training set that are also present in the test set, thus reducing the size of the training set from 210 music pieces down to 139 pieces. Following existing conventions [12, 18, 19, 41], all audio recordings are downsampled from 44.1 kHz to 16 kHz.

To demonstrate the effectiveness of our VAT model, we train our model using the following three versions of the MAPS dataset:

3.1.1 Full version. This version uses all 139 available pieces from MAPS as the labelled training set. To demonstrate the ability of leveraging unlabelled data using our VAT model, we use the training set from MAESTRO [19] as the unlabelled dataset (967 music recordings). The labelled training batch size N_l and the unlabelled training batch size N_{ul} are both 8.

3.1.2 Small version. In this version, only one folder (AkPnBcht, containing 23 non-overlapping songs) from MAPS is used as the labelled training set. We keep using the same 967 music recordings from MAESTRO as our unlabelled set for our VAT model. Again, N_l and N_{ul} are both 8 in this version.

3.1.3 One-shot version. Only one music recording (chp_op31 from the AkPnBcht folder) is used as the labelled training set. The unlabelled set consists of the same 967 music recordings from MAESTRO as the above two versions. Due to the fact that there is only one labelled training sample, N_l is 1 and the unlabelled training batch size N_{ul} remains 8.

3.2 MusicNet dataset

MusicNet [44] contains both audio recordings and annotations of various types of musical instruments such as those from the string family and the woodwind family. To prove that our model also works for different types of musical instruments, we perform our experiments on the following variations of MusicNet:

3.2.1 String version. In the official training set provided by MusicNet, there are 8 genres of music that contain string instruments. We select only one piece from each genre from the official training set, forming our own labelled training set. The remaining pieces of each genre are used as the unlabelled training set for our VAT framework. By doing so, there are eight labelled samples and 104 unlabelled samples in our training set. We pick four string pieces from the official test set provided by MusicNet as our test set. The IDs of the four pieces are 2191, 2628, 2106, and 2298 which are solo violin, accompanied violin, string quartet, and solo cello respectively. More details about the data splitting can be found in the supplementary material².

The labelled training batch size N_l and the unlabelled training batch size N_{ul} are both 8.

3.2.2 Woodwind version. Similar to the string version, we pick only one piece from six different woodwind genres from MusicNet as the labelled training set and use the remaining pieces in each genre as the unlabelled training set. This results in six labelled training samples and 21 unlabelled training samples. The official test set provided by MusicNet contains only two pieces (1819, 2416) from

the woodwind family, which belong to the Pairs Clarinet-Horn-Bassoon genre. We use these two pieces as our test set. Again, more details can be found in the supplementary material². N_l is 1 and N_{ul} is 8 in this version.

3.3 Data Processing

We extract Mel spectrograms on-the-fly from the audio clips using a GPU-based audio processing library nnAudio [11]. Following Hawthorne et al. [18], we use a Hann window size of 2,048, a hop size of 512, and 229 Mel bins as the parameters of our Mel spectrograms X_{spec} . To extract a fixed length spectrogram, we crop the audio clips into segments of 327,680 sample points using random sampling during each iteration, which results in Mel spectrogram with 640 timesteps, and 229 Mel frequency bins. We compress the magnitude of the spectrograms by taking the natural logarithm and then normalizing the magnitude for each spectrogram into the range $[0, 1]$. i.e. $X_{\text{spec}} \in [0, 1]^{640 \times 229}$.

As for our ground truth labels, we extract the onset, duration, and pitch information from the midi annotations to produce tsv files for the ground truth. These tsv files are read and converted into piano rolls in the form of a binary matrix $Y_{\text{label}} \in \{0, 1\}^{640, F}$. Since most musical instruments in the dataset are within the 88 notes range (note A0 to note C8), we use $F = 88$ in all our experiments.

3.4 Implementation Details

All models and experiments, including the baseline models, are implemented in PyTorch. To ensure transparency and fairness, we train all our models without tricks such as label smoothing [47], weighted cross entropy [18], and focal loss [29, 46]. We believe that these tricks would in general improve the transcription accuracy, and it is beyond the scope of this paper to explore this.

We adopt U-net models specifically designed for pitch detection [12, 21] and integrate them into the VAT framework [34]. While we follow mostly the same design as in [12], we modify the final layer of the decoder so that it has the flexibility to output two channels as shown in Figure 2. One of the channels is fed to a fully connected layer with sigmoid activation to predict the onsets $Y_{\text{onset}} \in [0, 1]^{T \times F}$, and the other channel is fed to a linear fully connected layer to obtain the features $Y_{\text{frame}} \in \mathbb{R}^{T \times F}$. The concatenated output $Y_{\text{onset}} \oplus Y_{\text{frame}}$ is fed to a relative local 1D self-attention layer [37, 40] to obtain the posteriorgram $Y_{\text{post}} \in [0, 1]^{T \times 88}$. We binarize the posteriorgram with a threshold of 0.5 to obtain the predicted piano roll $Y_{\text{roll}} \in \{0, 1\}^{T \times 88}$. If the two-channel output is used, we follow the inference method from the Onsets and Frames model [18] to obtain a refined piano roll by using both Y_{onset} and Y_{frame} to filter out notes that do not have an onset. Otherwise, we directly use the posteriorgram to obtain the piano roll. In addition, we also replace all of the LSTM layers in [12] with local relative self-attention layers, since it has been shown that self-attention layers perform as good as LSTM layers while providing the extra benefit of being able to train in parallel [45, 47].

We also modify the original VAT method [34] so that it works better for AMT. Firstly, since polyphonic AMT is a timestep-wise multiclass classification problem (multiple pitches can occur at the same time), we replace the Kullback–Leibler divergence (KL-div)

with binary cross entropy (BCE) when calculating the local distributional smoothness (LDS). Secondly, we normalise the adversarial vector r_{adv} along the timestep dimension as shown in Equation (2):

$$r_{\text{adv}} = \epsilon \left[\frac{g_1}{\|g_1\|_2}, \frac{g_2}{\|g_2\|_2}, \dots, \frac{g_T}{\|g_T\|_2} \right] \quad (2)$$

where ϵ is a parameter that controls the magnitude of the adversarial vector r_{adv} , and g_t for $1 \leq t \leq T$ is the timestep-wise gradient obtained from Equation (3)

$$g = \nabla_r \sum_i \text{BCE} \left[p(Y_i | X_{\text{spec}}, \hat{\theta}), p(Y_{\text{adv}} | X_{\text{spec}} + r, \hat{\theta}) \right]. \quad (3)$$

If the onsets prediction module is included, then $i = \{\text{onset, post}\}$. Otherwise, there is only one term in Equation (3), i.e. $i = \{\text{post}\}$. As in [34], the weight of the model is considered as a constant $\hat{\theta}$ when calculating the gradient g .

Once we obtain the adversarial vector r_{adv} , we can calculate the LDS. By the same logic as above, the LDS can contain either one or two terms depending on the model output:

$$\text{LDS}_* = \frac{\sum_i \text{BCE} \left[p(Y_i | X_{\text{spec}}^*, \hat{\theta}), p(Y_{\text{adv}} | X_{\text{spec}}^* + r_{\text{adv}}, \theta) \right]}{N_*}. \quad (4)$$

From Equation 4, we can see that the label Y_i^{label} is not required to calculate the LDS. Therefore, LDS is an unsupervised loss that can be calculated using both labelled spectrograms X_{spec}^l and unlabelled spectrograms X_{spec}^u . We will denote the LDS calculated using X_{spec}^l as LDS_l and the LDS calculated using X_{spec}^u as LDS_u . Unlike the original VAT [34], we normalise LDS_l and LDS_u by its respective batch size N_l and N_u , rather than summing both LDS_l and LDS_u together and normalize with $N_l + N_u$. By doing so, we prevent N_u from interfering with LDS_l and N_l from interfering with LDS_u .

3.5 Training Objective and Optimization

As mentioned in Section 3.4, we have the supervised objective L_l that requires labels, and the unsupervised objective L_u that does not require any label. The final objective L being minimized during training contains three terms as shown in Equation (7):

$$L_l = \sum_i \text{BCE} \left[Y_i, Y_i^{\text{label}} \right] + \sum_i \text{BCE} \left[Y_i^{\text{recon}}, Y_i^{\text{label}} \right] \quad (5)$$

$$L_u = \frac{\text{LDS}_l + \text{LDS}_u}{2} \quad (6)$$

$$L = L_l + \alpha L_u + L_{\text{recon}} \quad (7)$$

where α is the weighting for L_u , which is set to 1 throughout all our experiments; L_{recon} is the reconstruction loss mentioned in Section 2.2. We observe the same model behaviour as reported in [34], that is, controlling the ϵ in Equation (2) alone is sufficient to control the model performance without the need to change α .

To minimize the objective L , we use Adam [25] optimizer with a learning rate of 0.001 and a learning rate decay of 2% every 1,000 iterations. When training, our framework includes three forward passes during each iteration. One forward pass for L_l , one forward pass for LDS_l , and one forward pass for LDS_u . We define one epoch as 10 iterations. During the parameter search, we split our training set into 80% for training and 20% for validating. The optimal value for ϵ in Equation (2) is mostly within the range between 1 and 2,

and depends on the model architecture and the dataset. This value can be easily obtained after a few trials.

3.6 Evaluation Metrics

Following existing literature [12, 13, 18, 19, 23, 24], we report the frame-wise, note-wise, and note-with-offset-wise metrics to evaluate our model performance comprehensively. For note-wise metric, we use a onset tolerance of 50ms; for note-with-offset-wise metric, we use an offset tolerance of 50ms or 20% of the note duration, whichever is larger [2]. Readers are referred to Cheuk et al. [13] which explains the differences between these metrics in detail in their Section IV-C. In our experiments, we use the implementations from `mir_eval`¹ to calculate and report the above-mentioned metrics.

4 RESULTS

4.1 Effectiveness of VAT

We compare our proposed models to the Onsets and Frames model [18] and the Multi-Instrument AMT model [47] as they show good performance on the MAPS and MusicNet datasets respectively. We exclude the models proposed by Pedersoli et al. [36] and Thickstun et al. [43] in our results below since their performance is worse than the Multi-Instrument AMT model [47]. In Table 1-2, we use R to represent the reconstruction module, and O to represent the onset module. Therefore U-net-RO means that the U-net model contains both a reconstruction and onset module. The columns of the tables represent the precision (P), recall (R), and F1-score for each of the metrics mentioned in Section 3.6. Our proposed models and the baseline models are trained on the same labelled data, and only the proposed semi-supervised models are able to leverage the unlabelled data mentioned in Section 3.

4.1.1 Full MAPS. We can see that when using the VAT (row A3-A4, A7, A8), all three metrics generally improve compared to their respective counterparts without the VAT (row A1-A2, A5, A6). When using onset inference (A5-A8), the note-wise and note-with-offset-wise metrics are improved by at least 7 percentage points. The model using both the onset inference as well as our proposed framework (row A8) performs as good as the state-of-the-art Onsets and Frames model [18] (row 9) for this dataset.

4.1.2 Small MAPS. The middle part of Table 1 shows that when the number of labelled training samples is reduced by over 80% from 139 to 23 audio clips, the advantage of the VAT module becomes more obvious. Similar to the full MAPS dataset, the models with VAT module outperform their counterparts that do not use VAT. Moreover, our proposed framework (row B8) outperforms the Onsets and Frames model (B9) by 6, 5.1, 4.4 percentage points in terms of frame-wise, note-wise, and note-with-offset-wise F1-scores, which can be translated into improvements in performance of 11.5%, 8.1%, and 14.1% respectively.

4.1.3 One-shot MAPS. The bottom part of Table 1 shows that when we reduce the number of labelled training audio clip even further to only one, our proposed framework (C8) outperforms the Onsets and frames model (C9) by 23.7, 17, and 12.9 percentage points.

¹https://github.com/craffel/mir_eval

Between the models that use and do not use onset inference, we can see that onset inference has the tendency of decreasing the frame-wise F1-score while improving the note-wise F1-scores. This is due to the unreliability of the frame-wise metric [13, 18]. Cheuk et al. [13] has provided a few examples and shown that a high frame-wise score does not guarantee a good transcription. Nonetheless, these three experiments have shown that VAT is a very effective semi-supervised method, that allows the use of unlabelled training samples to greatly improve the model performance in cases where the number of labelled samples is scarce.

4.1.4 String MusicNet. The top section of Table 2 shows the performance of our proposed framework on the string subset of MusicNet (3.2.1). Although most models with the VAT outperform their counterparts without the VAT, U-net-RO VAT on row D4 performs worse than its counterpart without the VAT. Moreover, using the onset inference (row D1-4 and D9) does not improve the transcription accuracy in this setting, on the contrary, it worsens the model performance. There are two possible reasons for this. First, we believe that the onset inference only works well for piano only, and it cannot generalize well to other musical instruments such as those from the string and the woodwind family. Second, we believe that the onset labels for MusicNet are not completely accurate, since the labels are generated using dynamic time warping (DTW) [43]. Therefore, inaccurate onset labels might confuse the VAT. Using no labels might be better than using inaccurate labels, which is one of the advantages of using VAT.

Now, let us consider models (row D5-D8) that do not use onset inference. We will use the Multi-Instrument AMT model (row D10) [47], which is the state-of-the-art model for the MusicNet dataset at the time of writing, as the baseline model. Since the baseline model [47] is much deeper than the U-net model (row D5), it outperforms the U-net model. By applying the reconstruction module to the U-net model (row D6), the U-net model begins to outperform the baseline model. When we further apply VAT to the U-net models (row D7-D8), the transcription accuracy becomes even better. The best model, U-net-R VAT (row D8), outperforms the baseline model by 3.9, 11.1, and 11.3 percentage points in terms of frame-wise, note-wise, and note-with-offset-wise metrics.

4.1.5 Woodwind MusicNet. The bottom section of Table 2 shows the results for the woodwind subset of MusicNet (Section 3.2.2). Since the Onsets and Frames model does not work well for this dataset either, we did not spend time experimenting with it. Just like all of the results reported above, the VAT module is very effective in improving the transcription accuracy. The best model being the one with both the reconstruction and the VAT module (row E4), and it outperforms the baseline model by 10.3, 6.9 percentage points in terms of note-wise and note-with-offset-wise metrics. The improvement for the frame-wise metrics is not obvious, however, we must keep in mind that this is not a reliable metric to evaluate the transcription accuracy as pointed out previously in Section 3.6 as well as existing literature [13, 18].

4.2 Model Compactness

A comparison of number of trainable model parameters for the baseline models and the proposed models is shown in Figure 1. It can be

seen from the figure that a deep model does not necessarily yield a high transcription accuracy when the labelled training data is limited. The Onsets and Frames model [18] and the Prestack-Unet [36] have a high number of parameters, yet they do not perform well when the labelled data is scarce. While Thickstun's model [43] performs better than the two baseline models, its number of parameters is 10 times more than our proposed framework (U-net-R VAT). We use the Resnet-18 version of Prestack-Unet since the Resnet-32 version is too huge to run on our GPU. Another baseline model, the Multi-Instrument AMT model, performs better than the plain U-net model. With VAT, however, the U-net models already outperform the baseline model while keeping the number of trainable parameters low. We can also see that the VAT improves the model performance without adding extra parameters to the model. Therefore, VAT is a very effective method to improve the transcription accuracy by leveraging unlabelled training data when the labelled training data is limited.

5 APPLICATIONS

Our proposed semi-supervised framework allows for two important applications: continual learning and knowledge transfer to unseen music genres. We will discuss these two properties and their potential applications.

5.1 Continual Learning

The loss function of our proposed semi-supervised AMT framework contains a supervised term L_l and an unsupervised term L_{ul} . Even when we encounter new unseen, unlabelled data, we can still use this new data to minimize the unsupervised part of the model L_{ul} . That means, the proposed model can be *retrained* with any new data that was not collected before. Therefore, our proposed framework is capable of improving itself via new unlabelled data.

To confirm our framework's ability of continual learning, we take the string and woodwind subset of MusicNet as an example. We first train our models for 4,000 epochs (row 1 and 4 of Table 3, denoted as “4k”), and save the weights. These weights are then used as starting weights when we train the model for another 4,000 epochs with two different conditions: (1) without new data (row 2 and 5, denoted as “8k”); (2) using the test data as the unlabelled data as well as the existing data (row 3 and 6, denoted as “4k + 4k”). For the string subset of MusicNet, the model has already converged at 4,000 epochs, additional supervised training does not change the performance much. When we include the test data as the unlabelled data, it further pushes the accuracy around 1 percentage point higher. The same goes for the woodwind dataset. Although the improvement is relatively subtle at the moment, we plan to investigate ways to further improve this in future research. This property leads to the next application.

5.2 Case Study: Transcribing Unseen Genres

In some cases, we have some labels in one data domain, while the target domain we are interested in might not contain any labels at all. A model that can be trained on one domain and its knowledge then transferred to the target domain will be very useful. For example, we have some labelled data for classical woodwind music, but we want our model to be able to transcribe clarinet covers of Japanese pop

