

Timbre Separation of Monophonic Mixes Using Convolutional Neural Networks

Student: Alexander Albarosa

Supervisor: Claudio Balocco

Department of Engineering

University of Durham

Abstract—This paper demonstrates the capability of convolutional neural networks (CNN) to simultaneously separate and identify pitches of instruments in a mix via timbral discrimination. The separation of an electric guitar and a trumpet was investigated as their pitch range is overlapping and they have very different timbral characteristics. TensorFlow (TF) was used to construct and train CNNs on MIDI synthesised training datasets. Two labelling approaches were tested and the effect of a novel approach to the data used for training to emphasise timbral differences between instruments was analysed. The results suggest... TBD.

I. INTRODUCTION

Timbre separation is one of the key problems that must be solved to enable Automatic Music Transcription (AMT) for audio sources containing multiple instruments. AMT can be defined as "the design of computational algorithms to convert acoustic music signals into some form of music notation" [1]. AMT is a complex problem and has been tackled since the late 1970s, starting with the work by [2]; despite this, current solutions using machine learning are still significantly below the level of a human expert [3]; one particular challenge is that the combinatorially large space of possible musical sounds and notation necessitates vast, labelled datasets, however such datasets are scarce and creating them is a laborious task. The complexity associated with full AMT has led to its reduction into a set of sub-problems which individual pieces of research focus on; in [1] these are recognised as "(multi-)pitch estimation, onset and offset detection, instrument recognition, beat and rhythm tracking, interpretation of expressive timing and dynamics and score typesetting". Furthermore, there are several key factors that compound the difficulty of AMT; namely the: loss of information due to simultaneous polyphonic sound events from different sources (instruments) [4], ambiguity of musical notation, violation of statistical independence between source signals due to the metrical structure of music [1] and lack of annotated ground-truth transcriptions for supervised learning and testing (creating such annotations is very time consuming and requires high expertise [1]).

The most obvious use case for AMT is to significantly reduce the time taken for musicians to transcribe music (a frustration noted as early as 1743 by Jean Jacques Rousseau [5]); however, the International Society For Music Information Retrieval (ISMIR) research community has highlighted many other potential use cases of AMT such as: query-by-humming [6], automatic music tutoring [1], plagiarism detection [7]

and recreation of improvised performances [7,8]. There is a clear commercial opportunity for automatic music tutoring; *Yousician* [6] is a popular app that tailors the music tutoring process to an individual depending on their level. An appealing feature of *Yousician* to users is the ability to play, practice and receive feedback on songs featured in its library. AMT could significantly speed up the addition of songs into this library if it timbral separation capabilities such that the pitches of a particular instrument, such as a lead guitar, may be extracted from the song mix.

Recent systems for AMT can be broadly classified into two categories: new approaches using neural networks (NN) and traditional, unsupervised approaches not using NNs. State-of-the-art results have been obtained with both traditional methods and NNs in recent history using the Music Information Retrieval Exchange (MIREX) [7] evaluation metrics, however in the last few years much research has concentrated around NNs for their ability to learn high level features on their own [8], [9] and achieve state-of-the-art results [10].

AMT systems typically fall into one of the following four levels [1]:

- 1) Frame-level: the system outputs the fundamental frequencies present for each time frame for all time frames of the input audio.
- 2) Note-level: the system outputs frame-level pitch estimates connected over time to produce a piano-roll representation. Each note estimated has a pitch, onset time and offset time.
- 3) Stream-level: the system outputs note-level transcriptions for each instrument in the mixture by analysing the timbres present in the audio.
- 4) Notation-level: the system output is a human-readable score, this is the final goal of AMT.

This work falls mostly in the frame-level, however it also incorporates aspects of both the stream and note-levels: a piano-roll output was produced for each of the two instruments that were separated.

Recent frame-level approaches typically focus on polyphonic pitch estimation (a much harder task than monophonic pitch estimation [11]). Systems have varied in their methods using: traditional signal processing [12], [13], spectral factorisation techniques [14] (e.g. non-negative matrix factorisation (NMF) and NMF combined with probabilistic latent component analysis (PLCA) as seen in [15]), and NNs [16], [17]. According to [1], traditional signal processing methods are

faster but less accurate; [16] reports up to 9% improvement in pitch estimation accuracy using NNs over NMF and PLCA methods.

Note-level approaches are currently the main focus of AMT [11] with NNs achieving state-of-the-art results, especially when using convolutional neural networks (CNN) [3], [10], [18]. The work in [10] using CNNs beats out the best NMF approach in [19] by 1.5% and more traditional methods by approximately 30% in F-measure, defined in Section. III-C, benchmarks, showing the great promise of NNs (and CNNs) with regards to note-level transcription.

Little research has been done on either stream-level or notation-level AMT tasks. Stream-level AMT is closely related to audio source separation tasks (ASS) [20], and generally relies on recognising the timbre differences between instruments. ASS has been researched considerably though knowledge of it has yet to be widely applied to AMT tasks. The work in [21] is one of the few that has attempted stream-level AMT, though their approach uses NMFs and PLCA methods rather than NNs. One of the few approaches that utilises CNNs for a task approaching stream level transcription is [8], here a CNN is used for frame-level instrument prediction and for modelling the notes active in a frame; the task is considered as a multi-label classification problem. The input to the CNN consists of a CQT, a harmonic series feature (HSF) which represents the energy distribution of the harmonics of music notes (the authors believe that this is key in the perception of instrument timbre), and pitch information (they experiment with pitch information from ground-truth labels and from pitch estimations by the CNN in [10]). Another interesting feature is that they use 1D convolutions along time as opposed to 2D convolutions, which according to [22] may capture frequency and timbral information in each time frame better.

The approach taken in this report specifically focuses on the ability of CNNs to separate and identify pitches of instruments in a two-instrument mix. Two labelling approaches are tested and novel approach to the data used for training, detailed in Section. III-A, is analysed to see if timbral information can be successfully identified and used to extract an accurate piano-roll representation of an input mix.

II. THEORY

Following the success of CNN implementations in computer vision applications, e.g. blurring faces in Google StreetView [43] and in radiology to class tumour types [44], the MIR community has attempted to leverage similar methods for AMT. In order to enable the use of CNNs for AMT, audio is typically processed into spectrograms (a 2D time-frequency representation) and labelled. Where images (and image datasets) are readily available online to train computer vision networks, data for AMT applications is scarce at best. The rest of this section covers the fundamental theory required to understand how CNNs work and how they can be used for the separation of instruments by timbre.

A. Convolutional Neural Networks

A CNN is a neural network that uses convolutional layers to output feature maps using learned, weighted filters; each

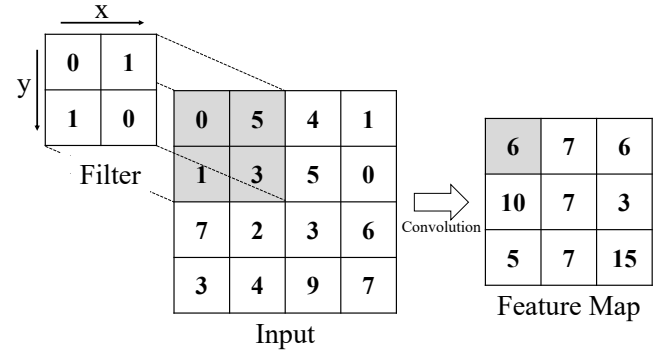


Fig. 1: Visual representation of the calculation of a feature map from a simple 2D filter and input using a stride of 1 in the x and y directions.

map represents a feature extracted at all locations in the input [23]. Convolutional layers are typically used in conjunction with pooling layers (to save computation resources) and this unit may be repeated many times before the outputs are fed into fully connected layers to allow for classification. In supervised training (where the network is trained using labelled inputs): stochastic gradient descent (SGD) is used to minimise the difference between the desired and actual output of the network; the filters in each convolutional layer are updated simultaneously by the learning procedure and back-propagation is used to compute the gradients [23]. For a mathematically rigorous explanation regarding the extraction of feature maps from the input and the learning procedure of CNNs the reader is directed to [23]–[25], although a brief explanation of the convolution and back-propagation mechanisms are provided here. The operation of convolution between an input and a filter can be seen in Fig. 1; the highlighted number on the feature map is simply the sum of the element-wise multiplication of the filter with the highlighted section of the input, the window. Using a stride (the distance which the filter ‘hops’ from one convolution to the next), $s = (1, 1)$ in the x and y direction respectively, the filter first moves horizontally across the input, hopping one element further along each convolution, once the filter reaches the end of the input row it then, like a type-writer, returns to the first column displaced by one row in the y -direction and repeats the process.

The result is a feature map of size $(3, 3)$. More generally, for an input of size (i_x, i_y) and a filter of size (k_x, k_y) , using a stride (s_x, s_y) the output is a feature map of size $((i_x - k_x)/s_x + 1, (i_y - k_y)/s_y + 1)$.

The second key idea behind CNNs and neural networks in general is back-propagation. Back-propagation is the enabling idea behind SGD that allows networks to converge to local minima of the cost function, C , based on the negative gradient of C , or $-\nabla C$. Every time a new batch of outputs is produced by the network, the difference between the produced output and the desired output is calculated using the C (the particular C used depends on the nature of the task). The goal of the network is then to minimise C by adjusting the weights and biases which calculate the activation of each neuron

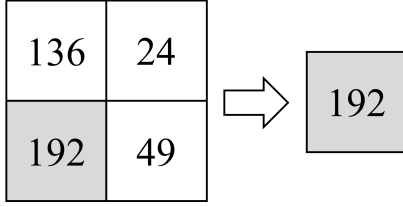


Fig. 2: The operation of a 2D max-pooling layer to reduce the size of the input.

corresponding to the outputs. The magnitude by which each weight and bias is changed should be changed is determined by $-\nabla C$; this quantity is calculated using back-propagation, which is achieved using a specific case of reverse automatic-differentiation to compute the partial derivatives required to calculate the sensitivity of the output neurons to any given weight or bias [24], [26]. These values allow the network to tune its weights and biases, in a controlled way, to minimise C and thus produce the desired output.

There are certain different training methods in which back-propagation is used, a typical supervised learning approach may use a dense layer as the final layer to classify (or discriminate) the input sample then, using C , all the weights in the layers in the network are tuned accordingly; however other methods, such as adversarial, autoencoder and reinforcement learning can also be used. In adversarial learning two separate NN are trained, a generator and a discriminator; the discriminator learns to recognise the difference between the output from the generator and that of a real input that the generator is attempting to mimic, the task of the generator is to generate a "fake" output from input noise such that the discriminator is tricked into thinking that the "fake" input is in fact real at least 50% of the time [27]. Such networks have found use in music synthesis [28] and can be used to improve transcription networks [29].

One of the main advantages of using CNN is that they allow for progressively smaller features to be extracted as more convolutional layers are used, this increases the granularity of the information available for classification of the input. However the use of more layers requires increased computational resources; to offset this effect pooling layers are used to reduce the size of the feature maps between convolutional layers. The operation of a common pooling layer, *MaxPooling2D* [30], can be seen in Fig.2.

A few limitations of CNN are the requirement for large, accurately labelled datasets and the tendency to overfit the training data, though overfitting can be mitigated by the use of dropout layers which force the network to generalise better [17].

In the context of timbre separation and AMT, spectrograms (2D frequency-time representations of audio) have been used as input to CNNs in numerous studies operating at the frame-level of AMT [10], [11], [17]. Spectrograms are in their nature inherently suited to be used as input to CNNs: music events are highly correlated over both time and frequency [1] and thus using spectrograms, CNNs can learn high-level features and relationships temporally and in frequency [10]. Specifically

regarding instrument separation by timbre and more generally frame-level AMT; it is thought that CNNs can be used to learn the unique energy distributions of partials for different instruments [8] from spectrograms. Unfortunately, there have not been many studies into frame-level AMT and thus there is a palpable lack of data available for the task [31]. One way around this is to synthesise data using MIDI instruments. Though MIDI instruments will never be able to sound exactly like their real counterparts, they allow for simple generation of audio data through which easy labelling is facilitated by *Python* packages such as *pretty_midi* [32] which make parsing MIDI files into labels very easy. This ease of labelling is extremely important for AMT where making sure that labels are accurately aligned to spectrograms is a key, and time consuming issue to ensure efficient training of CNNs [33]. Another advantage of synthesising and labelling one's own data for the task of timbre separation is the ability to accurately control the balance of data for different instruments in the training dataset.

It is important to remember that MIDI synthesised data can only go so far: there can be large differences in timbre between an instrument synthesised digitally using MIDI sound fonts and the same instrument recorded with a microphone in a studio. To enable NN to generalise well to real-world recordings, real-instrument data should be used for training. If such data is unavailable then the network should be tested on real-instrument recordings to gain valuable insight as to the applicability of the network in real-instrument scenarios.

In this study it is proposed that in order to best train a neural network to be able to separate instruments by timbre and estimate their pitch, there are two key considerations to be made with regards to the training data. The data regarding each instrument must include: the pitch in frames where it is playing, a label symbolising a rest in all frames where it is not present, including data concerning all other instruments that could co-exist in a mix with the given instrument to be separated. The latter consideration is proposed as a way of clearly communicating to the network that one instrument is not another instrument thus forcing it to recognise the timbral differences between instruments, leading to improved separation. With regards to the task of pitch estimation, the problem is considered as a multi-label classification problem with the loss function:

$$C_{cce} = - \sum_{c=1}^n b_{o,c} \log(p_{o,c}) \quad (1)$$

where c represents the class, n is the number of possible classes, b is a binary indicator as to whether the class classification c is correct for output o and p is the predicted probability that o belongs to c [34]. Treating pitch estimation in this way works well for Western styles of music as each class can be defined as corresponding to a specific MIDI pitch from 0 to 127 inclusive, representing the octave range C0 to G9.

B. Spectrograms

Seen as spectrograms are the key data enabling the use of CNNs for AMT tasks, it is important to understand how they

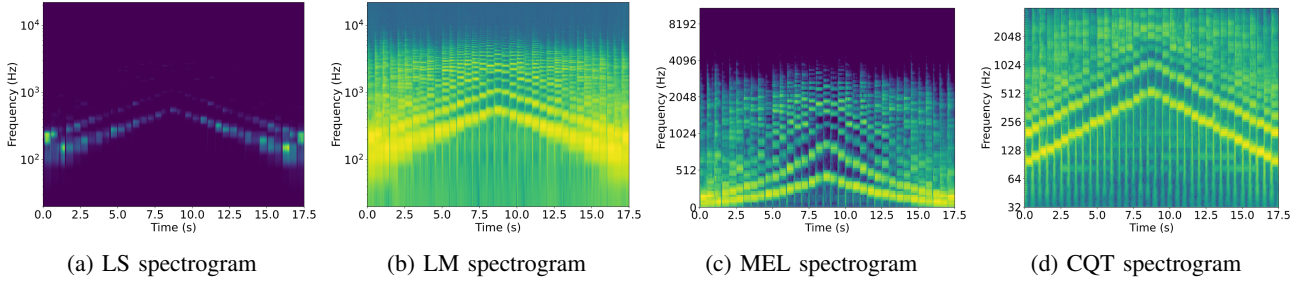


Fig. 3: The four commonly found spectrogram representations in AMT systems all displaying an ascending and descending G major scale played on an electric guitar.

are constructed and the good practices behind their creation. As mentioned in the Section. II-A spectrograms are 2-D representations of audio that can be labelled and used as input to a CNN for training. They are produced using the short-time Fourier transform and there are four common ways to scale the result as a graph of frequency against time:

- 1) Logarithmically spaced frequency bins (LS).
- 2) Logarithmically spaced frequency bins and logarithmically scaled magnitude (LM).
- 3) Mel-scaled frequency bins (MEL).
- 4) Constant Q-Transform (CQT).

Fig. 3 shows these four types of spectrograms. Each one was computed at a sample rate of 44.1 kHz with an STFT window length of 1024 samples (giving a frequency bin resolution of 43 Hz) and hop size of 512. Fig.3a and Fig.3b were manually implemented in python with help from the `pyFFTW` package [35], a wrapper around the popular C library `FFTW` that implements fast Fourier transform algorithms. Fig.3c and Fig.3d were computed with aid from the `librosa` package [36], a package specifically designed for music and audio analysis for music information retrieval.

Often research papers focus on one type of spectrogram, although some studies, such as [17] and [11], show comparisons of results using the different spectrograms for training and testing. In practice, there is little difference between the spectrogram types, all apply some kind of logarithm scaling to the frequency axis, and there seems to be no standard or rule as to which representation to use, though [37] suggests that a CQT spectrogram is better suited for music tasks as the frequency axis is linear in pitch (this also allows for easier data augmentation with regards to pitch shifting samples). Where differences can be found in research lie in the window size chosen for the STFT: there is an inverse relationship between time and frequency resolution, i.e. shorter windows (increasing temporal resolution) lead to larger frequency bins (decreasing frequency resolution) and vice versa, thus there is a balance to be struck for AMT systems to be able to extract sufficient timing information for note values and sufficient frequency information for pitch detection. For AMT applications it is common [1] to aim for a temporal resolution (hop length) in the order of 10ms for a given sampling frequency, this is achieved by choosing a small hop length between frames, however this can lead to large overlap which can cause ambiguity. To reduce computational overhead it

is also common to downsample audio files [10], [14], [16], [38], [39], this does however come at the cost of reducing the highest frequency that can be recovered from the signal (due to Nyquist's sampling theorem).

III. METHOD

In order to test the hypotheses made in Section. II-A, it was necessary to create tailored training and testing data. To enable easier alignment of labels and spectrogram frames the decision was made to use 120BPM audio sources only. Furthermore, to further simplify training and testing, only monophonic, two instrument mixes were considered; the main focus being a mix of trumpet and electric guitar. This combination is good to test for the task of timbre separation due to the pitch range overlap of the instruments and their very different timbral characteristics. The electric guitar features in all mixes as real guitar data could be recorded and tested using the *Native Instruments* Komplete Audio 1 USB Audio Interface. The neural networks were created with the widely used TensorFlow (TF) [40] package.

A. Dataset Creation and Data Preprocessing

The process of training dataset creation was broken down into five major steps:

- 1) Creation of MIDI 'instrument ascension' (see text for definition).
- 2) Exporting MIDI file as .wav file.
- 3) Creation of spectrograms.
- 4) Spectrogram labelling.
- 5) Dataset creation.

The first step involved the use of `pretty_midi` to rapidly produce an 'instrument ascension' in the format of a .mid file. An instrument ascension was considered to be a sequence of crotchets starting from the lowest, and ending at the highest, note in an instrument's range. In the case of two instrument mixes, the instrument ascension consisted of every monophonic note combination between the two instruments. The .mid instrument ascension was then opened in *Musescore 3* (a popular music-score editing software). Soundfonts (a mixture of default and third party [41], [42]) were applied to the ascensions, then using *Musescore 3* the MIDI files were exported as WAV files (sampled at 44.1kHz). Using *Python* and the packages: `librosa`, `pyfftw` and `matplotlib`, the WAV files were down-sampled to 32.768kHz, spliced

into sections containing one note at a time and converted into LM spectrograms covering frequencies from 20Hz to 16kHz saved as .png images of size 128x512x3. This relatively large image size was chosen so that as much spectrogram information as possible was available to the NN, furthermore the sampling frequency was kept relatively high to retain the energy distribution of harmonics thought to be key to perceiving instrument timbre [8]. With regards to labelling the pitch of the instrument, two distinct approaches were taken: one used a "dual-label" containing two one-hot encoded, size 129 bit arrays, corresponding to the two separate instruments, and the other simply used one, similarly encoded, size 129 bit array to label the pitch of the desired instrument, "single-label"; Fig. 4 shows the different labelling approaches. The keen-eyed reader may have noticed that this is one bit longer than the number of possible MIDI pitches; this extra bit was used to signify the lack of presence of the instrument in the spectrogram. Once the spectrograms were assigned labels they stored in the format of a HDF5 dataset using the `hdf5` [43] package. The HDF5 format was chosen due to its superior performance when reading large amounts of images [44]. Special care was taken to ensure that datasets contained approximately (within 10% error) the same number of data points for each single instrument and combinations of instruments so as to ensure accurate models.

The data used for testing underwent the same steps, except for the first. The MIDI files used for the creation of the testing dataset are from a specific subset of the Lakh-MIDI dataset [45], [46] known as the "clean MIDI subset". This subset contains the music from songs of different genres in MIDI format. Eight of these MIDI files were adapted to include just two monophonic melodies, corresponding to the instruments that the network was trained to separate, and set to a tempo of 120BPM before being exported as WAV files; furthermore the audio of each of the instruments playing alone was exported in the same way to allow for more rigorous analysis of the performance of the networks. As a result there were three testing sets, each containing 35 minutes of music, see Table I. There was a conscious effort made to use varying soundfonts when compared to the training data, and in the case of melodies with guitars, some were rendered using a clean guitar soundfont and then processed using *GuitarRig 5* in *Ableton Live 10 Lite* to simulate a more realistic sound sample.

TABLE I: The three datasets used for testing.

Dataset	Instrument Mix	
	Guitar	Trumpet
Combined-Test	✓	✓
Guitar-Alone-Test	✓	
Trumpet-Alone-Test		✓

B. Training Process

The training process for every tested neural network configuration consisted of a K-fold cross validation with $K = 5$ (implemented using the `scikit` [47] package), with each fold running through 20 epochs using a batch-size of 32 (where

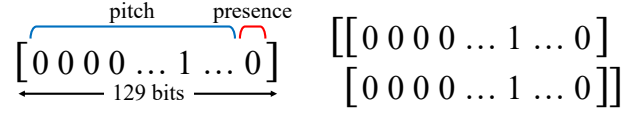


Fig. 4: The labelling methods used. The "single-label" approach can be seen on the left along with the bits allocated for the pitch and whether or not the instrument was present (0 means that the instrument is present). The "dual-label" approach can be seen on the right and simply contains two copies of the "single-label" format, each corresponding to one instrument in the mix.

memory constraints allowed). K-fold cross validation refers to the process of first shuffling and then splitting the training dataset into K distinct folds each with a training:validation split of $(1-1/K:1/K)$; when $K = 5$ this corresponds to an 80:20 split. This is done to prevent "lucky" data producing an optimistic accuracy; the accuracy over the K-folds is averaged to get a more realistic accuracy and for the purposes of this research, the fold with the best training accuracy was used to assess performance on the testing dataset. Unfortunately this process is time consuming due to the increased number of epochs required, thus to enable faster training, GPU acceleration in TF, with an NVIDIA GTX 1050 GPU with 4GB of video memory, was used.

C. Metrics

The performance of each network was analysed via direct comparison of the piano-roll produced from the raw MIDI file versus the piano-roll extracted from the predictions of the NN (facilitated by the known duration of the spectrograms in time). A piano-roll is pitch-time representation of audio data at a certain sampling frequency (this was chosen to be 100Hz, corresponding to a time resolution of 10ms, five times stricter than the 50ms used in the MIREX evaluation metrics [7]). Figure ?? shows an example of a piano-roll.

The accuracy of the predicted piano-roll was measured by four activation metrics: "correct" (in time and pitch), "wrong" (correct in time but not in pitch), "false" (a note was found when there should have been none) and "missed" (there was no note found when there should have been one). From these, the precision (P), recall (R) and F-measure (F) [48]:

$$P = \frac{correct}{correct + wrong + false} \quad (2)$$

$$R = \frac{correct}{correct + missed} \quad (3)$$

$$F = 2 \frac{PR}{P + R} \quad (4)$$

were calculated to give a more holistic view of the performance of each network, considering that wrong, false and missed notes are highly undesirable in musical scenarios [1]. For every CNN, these metrics were calculated across the three testing datasets (as described in Section III-C and the results were averaged to obtain the overall performance of any one CNN.

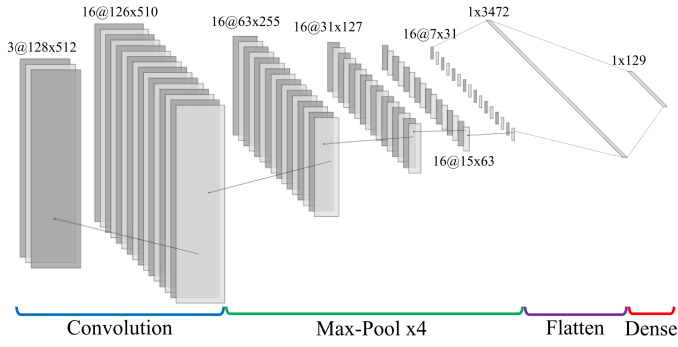


Fig. 5: Architecture of the shallow network used as a baseline for performance. The input is the 128x512x3 spectrogram; the convolutional layer produces 16 feature maps; the 2D max-pooling layers use 2x2 filters to reduce the size of the input to 7x31x16 before flattening; the dense layer is fully connected and has 129 output neurons.

TABLE II: The datasets used for training

Dataset	Frame Size	Number of Samples			
		Guitar	Trumpet	Combined	Total
Guitar Alone	8192	742	-	-	742
Trumpet Alone	8192	-	225	-	225
Both Alone	8192	742	225	-	967
Combined_1	8192	742	225	3243	4210
Combined_2.0	8192	2968	3150	3243	9361
Combined_2.1	4096	5936	6300	6486	18722

D. Training and CNN Configurations

Initially, three networks were trained and tested on three different datasets each. The baseline neural network used with the single label approach was considered to be a shallow network with the structure seen in Fig. 5. The output layer is a fully-connected dense layer, using a *softmax* activation with 129 outputs. The *softmax* activation ensures that the probabilities of all the output neurons sum to 1, this enabled easy extraction of the predicted MIDI pitch of the instrument using the `numpy` [49] package to find the index of the largest probability in the output pitch array. An analogous two-label network was created to serve the same purpose as a baseline network. The two other networks had one and three additional convolutional layers before flattening. The convolution filter size used was 3x3 for all convolutional layers. The spectrograms in the datasets used at this stage had a frame size of 8192 samples (corresponding to a duration of 0.25s) using a hop-length of 512 samples. The datasets used for training are summarised in Table II.

The best performing dataset (within the first four rows) was taken forward for both the label configurations; the dataset was then balanced resulting in an equal distribution of solo and combined instrument data, see row 5 of Table II. This dataset was then tested again on the three initial NN configurations and the best network and label configuration was noted.

Following this the effect of using spectrograms with hop-lengths of 32 and 1024, and frame sizes of 4096, 2048 (corresponding to a duration of 0.25s and 0.125s respectively), were tested on the best network and label configuration to determine whether the hop-length and window used had an

effect on the ability of the network to separate by timbre and distinguish pitch. This allowed for the determination of the best representation of data for further experiments with varying NN configurations.

Once the baseline had been trained, experiments concerning the architecture of the networks, and filter size of the convolutional layers were also conducted. If a certain new configuration performed worse than any previous configuration it was discarded.

IV. RESULTS

Table III shows the performance of the best network and training dataset on the testing dataset, as well as their training accuracy. Fig. 6 shows the structure of the best network. Include figures showing a desired piano roll output and the

TABLE III: Performance of the best networks for the single and dual label approaches.

Dataset	Network	Training Acc.	Testing Acc.

network predicted piano roll output.

A. Effect of Input and Training Data

Throughout the testing the significance of the right combination and format of the input training data was apparent. The initial networks were tested on the datasets present in the first five rows of Table II, the results can be seen in Table IV. The shallow network performed best when trained on the balanced "Combined_2.0" dataset, however it should be noted that it suffered from reduced guitar precision compared to the "Both Alone" dataset, potentially meaning that the CNN has greater difficulty distinguishing whether the guitar is in fact alone or playing at the same time as the trumpet, than it does under the same conditions for a trumpet. However the results here are promising as they suggest that the CNN can recognise the difference between the timbres of the guitar and the trumpet when supervised with data for all possible monophonic pitch combinations. It should also be noted that the shallow network performed the best out of the the initial networks beating out those with more convolutional layers.

TABLE IV: Performance of the shallow network in Fig. 5 on the testing datasets for the different training datasets all using spectrograms with a hop length of 512 samples.

Dataset	Guitar		Trumpet		Average	
	P	R	P	R	F	Accuracy
Guitar Alone	0.529	1.000	0.333	0.333	0.509	0.431
Trumpet Alone	0.333	0.333	0.436	0.999	0.454	0.385
Both Alone	0.804	0.870	0.879	0.855	0.843	0.749
Combined_1	0.703	0.990	0.875	0.998	0.877	0.786
Combined_2.0	0.746	0.991	0.889	0.995	0.894	0.813

The shallow network was then tested with a variety of spectrogram parameters as described in Section III-D. It was found that spectrograms with a frame size and hop length of 4096 and 512 samples performed best. Table V shows that the average accuracy was improved by approximately 5% by reducing the frame size from 8192 to 4096 samples, with the

TABLE V: The effect of different frame sizes and hop lengths on the performance of the "dual-label" shallow network using the Combined_2 dataset type.

Frame Size	Hop Length	Guitar		Trumpet		Average	
		P	R	P	R	F	Acc.
8192	32	0.772	0.984	0.894	0.994	0.902	0.826
8192	512	0.746	0.991	0.889	0.995	0.894	0.813
8192	1024	0.754	0.984	0.889	0.994	0.895	0.816
4096	32	0.810	0.989	0.908	0.993	0.919	0.853
4096	512	0.836	0.988	0.912	0.993	0.928	0.868
4096	1024	0.836	0.988	0.898	0.997	0.925	0.862

gains of 10% seen in the Combined-Test and Guitar-Alone-Test datasets, however it is likely that the increase in accuracy is as a result of the increased time resolution obtained by decreasing frame size as opposed to the network becoming better at distinguishing timbre. This conclusion is supported by the fact that the percentage of "false" activations remained largely unchanged, in fact the gains in accuracy came from the reduction in the percentage of "wrong" activations; though this illustrates the importance of high temporal resolution in AMT.

Somewhat surprisingly the number of "false" activations did not decrease when the hop-length was reduced to 32 samples; the increased overlap was hoped to increase the resolution of timbral information available in the spectrogram by showing a more accurate decay of the power distribution of the note harmonics in time; however the failure to capitulate on this extra information may have been down to the use of a shallow network, that was incapable of extracting meaningful relations from the data.

B. Effect of Network Configurations

The results of testing on the three initial networks (described in Section III-D) demonstrated that while deeper CNN learned faster (which is expected behaviour), they performed slightly worse overall than the shallow network in Fig. 5. The biggest performance drop came from the Trumpet-Alone-Test dataset, where the shallow network outperformed the deeper networks by approximately 10-15% in F results averaged across both the instruments, the deeper networks especially struggled to infer that no guitar was present at all in these recordings; suggesting that the deeper CNNs struggled to learn the differences between the timbres of the guitar and trumpet. To counteract this, networks were created with convolutional layers with long, thin filters, one along time and the next across frequency (and vice versa), to attempt better capture the timbral characteristics of the instruments; this choice was motivated by the 1-D convolutions used in [22] and the rectangular filters used in [50].

C. Effect of Labelling Approach

The two labelling approaches were found to differ minimally (<1%) in accuracy in all the initial tests. As a result, the "dual-label" approach was used thereafter as it allowed for joint prediction of the guitar and trumpet by one network, reducing the amount of processing required both for dataset

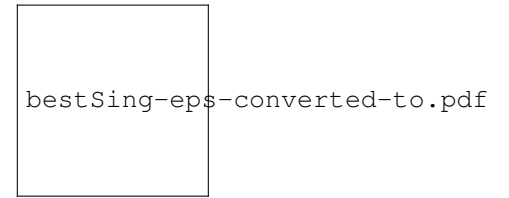


Fig. 6: Structure of the best network.

creation and for piano-roll extraction. However the "single-label" approach is more versatile: a network could in theory be trained to pick one instrument out of a mix of instruments if a "Combined" type dataset were used; this could allow for a user to specifically select the instruments that they would like to be transcribed out of a mix of instruments.

V. CONCLUSION

This networks and datasets analysed in this paper demonstrate the ability of CNNs to achieve timbre separation of instruments in mixes by learning timbral features present in spectrograms. An accuracy of X is achieved showing the potential of the proposed CNN topologies with labelled spectrograms. Furthermore a coding framework has been developed to quickly generate MIDI data and the training datasets created and used have been made publicly available. There are some obvious limitations to this approach, the limitation of tempo to 120BPM and lack of applicability to polyphonic sources, however it is hoped that this approach could be combined with ideas from other areas of AMT research to make progress towards achieving AMT for multiple instrument sources. Other avenues for further work could involve creating accurately aligned spectrogram data using real instruments; using an NN model to extract pitch information in real-time; using a similar approach to the one proposed here to separate an instrument of interest from mixes it is commonly present in (such as a trumpet in a jazz band) or even separating two of the same instruments from each other. Nonetheless this work serves as an encouraging step towards being able to achieve stream-level AMT and provides a simple way to produce accurately aligned MIDI generated data for machine learning.

REFERENCES

- [1] E. Benetos, S. Dixon, Z. Duan, and S. Ewert, "Automatic music transcription: An overview," *IEEE Signal Processing Magazine*, vol. 36, no. 1, pp. 20–30, 2019.
- [2] M. Piszczalski and B. A. Galler, "Automatic music transcription," *Computer Music Journal*, vol. 1, no. 4, pp. 24–31, 1977.
- [3] M. A. Román, A. Pertusa, and J. Calvo-Zaragoza, "An End-to-end Framework for Audio-to-Score Music Transcription on Monophonic Excerpts," in *Proceedings of the 19th International Society for Music Information Retrieval Conference*, (Paris, France), pp. 34–41, ISMIR, Sept. 2018.
- [4] M. Roman, A. Pertusa, and J. Calvo-Zaragoza, "A Holistic Approach to Polyphonic Music Transcription with Neural Networks," in *Proceedings of the 20th International Society for Music Information Retrieval Conference*, (Delft, The Netherlands), pp. 731–737, ISMIR, Nov. 2019.
- [5] J. Rousseau, *Dissertation sur la musique moderne*. HACHETTE LIVRE, 1743.
- [6] "Yousician." <https://www.yousician.com>. (accessed Apr. 19, 2021).
- [7] J. S. Downie, "Mirex home." https://www.music-ir.org/mirex/wiki/MIREX_HOME. (accessed Apr. 19, 2021).

- [8] Y.-N. Hung and Y.-H. Yang, "Frame-level Instrument Recognition by Timbre and Pitch," in *Proceedings of the 19th International Society for Music Information Retrieval Conference*, (Paris, France), pp. 135–142, ISMIR, Sept. 2018.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- [10] C. Hawthorne, E. Elsen, J. Song, A. Roberts, I. Simon, C. Raffel, J. 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, ISMIR 2018, Paris, France, 2018*, 2018.
- [11] M. A. Román, A. Pertusa, and J. Calvo-Zaragoza, "Data representations for audio-to-score monophonic music transcription," *Expert Systems with Applications*, vol. 162, p. 113769, 2020.
- [12] V. Emiya, R. Badeau, and B. David, "Multipitch estimation of piano sounds using a new probabilistic spectral smoothness principle," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 6, pp. 1643–1654, 2009.
- [13] L. Su and Y.-H. Yang, "Combining spectral and temporal representations for multipitch estimation of polyphonic music," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 10, pp. 1600–1612, 2015.
- [14] F. Pedersoli, G. Tzanetakis, and K. M. Yi, "Improving music transcription by pre-stacking a u-net," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 506–510, 2020.
- [15] P. Smaragdis, B. Raj, and M. Shashanka, "A probabilistic latent variable model for acoustic modeling," in *In Workshop on Advances in Models for Acoustic Processing at NIPS*, 2006.
- [16] S. Sigita, E. Benetos, and S. Dixon, "An end-to-end neural network for polyphonic piano music transcription," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 5, pp. 927–939, 2016.
- [17] R. Kelz, M. Dorfer, F. Korzeniewski, S. Böck, A. Arzt, and G. Widmer, "On the potential of simple framewise approaches to piano transcription," in *Proceedings of the 17th International Society for Music Information Retrieval Conference, ISMIR 2016, New York City, United States, August 7-11, 2016* (M. I. Mandel, J. Devaney, D. Turnbull, and G. Tzanetakis, eds.), pp. 475–481, 2016.
- [18] C. Thomé and S. Ahlback, "Polyphonic pitch detection with convolutional recurrent neural networks," in *18th International Society for Music Information Retrieval Conference, Suzhou, China, 2017*.
- [19] L. Gao, L. Su, Y. Yang, and T. Lee, "Polyphonic piano note transcription with non-negative matrix factorization of differential spectrogram," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 291–295, 2017.
- [20] E. Manilow, P. Seetharaman, and B. Pardo, "Simultaneous separation and transcription of mixtures with multiple polyphonic and percussive instruments," *CoRR*, vol. abs/1910.12621, 2019.
- [21] V. Arora and L. Behera, "Multiple f0 estimation and source clustering of polyphonic music audio using plca and hmrf," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 2, pp. 278–287, 2015.
- [22] S.-Y. Chou, J.-S. Jang, and Y.-H. Yang, "Learning to recognize transient sound events using attentional supervision," in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pp. 3336–3342, International Joint Conferences on Artificial Intelligence Organization, 7 2018.
- [23] Y. LeCun, K. Kavukcuoglu, and C. Farabet, "Convolutional networks and applications in vision," in *Proceedings of 2010 IEEE International Symposium on Circuits and Systems*, pp. 253–256, 2010.
- [24] A. G. Baydin, B. A. Pearlmutter, and A. A. Radul, "Automatic differentiation in machine learning: a survey," *CoRR*, vol. abs/1502.05767, 2015.
- [25] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, p. 2278–2324, 1998.
- [26] A. G. Baydin, B. A. Pearlmutter, and A. A. Radul, "Automatic differentiation in machine learning: a survey," *CoRR*, vol. abs/1502.05767, 2015.
- [27] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," 2014.
- [28] C. Hawthorne, A. Stasyuk, A. Roberts, I. Simon, C.-Z. A. Huang, S. Dieleman, E. Elsen, J. Engel, and D. Eck, "Enabling factorized piano music modeling and generation with the maestro dataset," 2019.
- [29] J. W. Kim and J. P. Bello, "Adversarial learning for improved onsets and frames music transcription," in *20th International Society for Music Information Retrieval Conference, Delft, The Netherlands, 2019*.
- [30] Keras, "Keras documentation: Maxpooling2d layer." https://keras.io/api/layers/pooling_layers/max_pooling2d/. (accessed: Apr. 19, 2021).
- [31] Y. Hung, Y. Chen, and Y. Yang, "Multitask learning for frame-level instrument recognition," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 381–385, 2019.
- [32] C. Raffel and D. P. W. Ellis, "Intuitive analysis, creation and manipulation of midi data with pretty_midi," in *Proceedings of the 15th International Conference on Music Information Retrieval Late Breaking and Demo Papers, 2014*.
- [33] J. Thickstun, Z. Harchaoui, D. P. Foster, and S. M. Kakade, "Invariances and data augmentation for supervised music transcription," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2241–2245, 2018.
- [34] "ML glossary documentation: Loss functions." https://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html. (accessed: Apr. 19, 2021).
- [35] H. Gomersall, "pyfftw." <https://github.com/pyFFTW/pyFFTW>. (accessed: Apr. 19, 2021).
- [36] Brian McFee, Colin Raffel, Dawen Liang, Daniel P.W. Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto, "librosa: Audio and Music Signal Analysis in Python," in *Proceedings of the 14th Python in Science Conference* (Kathryn Huff and James Bergstra, eds.), pp. 18 – 24, 2015.
- [37] J. C. Brown, "Calculation of a constant q spectral transform," *The Journal of the Acoustical Society of America*, vol. 89, no. 1, pp. 425–434, 1991.
- [38] C. Chafe, D. A. Jaffe, K. Kashima, B. Mont-Reynaud, and J. O. Smith, "Techniques for note identification in polyphonic music," in *Proceedings of the 1985 International Computer Music Conference, Burnaby, B.C., Canada*, International Computer Music Association, International Computer Music Association, 1985.
- [39] E. J. Humphrey and J. P. Bello, "Rethinking automatic chord recognition with convolutional neural networks," in *2012 11th International Conference on Machine Learning and Applications*, vol. 2, pp. 357–362, 2012.
- [40] M. Abadi, A. Agarwal, and P. Barham, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015. Software available from tensorflow.org.
- [41] Z. Demircan, "Flamestudios guitar soundfonts in sf2 format." <https://www.hedsound.com/2019/07/flamestudios-guitar-soundfonts-in-sf2.html>. (accessed: Apr. 20, 2021).
- [42] S. Papel, "Papelmedia final sf2 xx1 - trumpet." <https://www.polyphone-soundfonts.com/documents/20-brass/199-papelmedia-final-sf2-xx1-trumpet>. (accessed: Apr. 20, 2021).
- [43] "The hdf5 file library and format." <https://www.hdfgroup.org/solutions/hdf5/>. (accessed: Apr. 19, 2021).
- [44] S.-H. Lim, S. R. Young, and R. M. Patton, "An analysis of image storage systems for scalable training of deep neural networks," 4 2016.
- [45] C. Raffel, "The lakh midi dataset v0.1." <https://colinraffel.com/projects/lmd/>. (accessed: Apr. 20, 2021).
- [46] C. Raffel, *Learning-Based Methods for Comparing Sequences, with Applications to Audio-to-MIDI Alignment and Matching*. PhD thesis, Columbia University, 2016.
- [47] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, P. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [48] M. Sokolova, N. Japkowicz, and S. Szpakowicz, "Beyond accuracy, f-score and roc: A family of discriminant measures for performance evaluation," in *AI 2006: Advances in Artificial Intelligence* (A. Sattar and B.-h. Kang, eds.), (Berlin, Heidelberg), pp. 1015–1021, Springer Berlin Heidelberg, 2006.
- [49] C. R. Harris, K. J. Millman, S. J. van der Walt, and R. Gommers, "Array programming with NumPy," *Nature*, vol. 585, pp. 357–362, Sept. 2020.
- [50] J. Pons, O. Slizovskaia, R. Gong, E. Gómez, and X. Serra, "Timbre analysis of music audio signals with convolutional neural networks," in *25th European Signal Processing Conference (EUSIPCO)*, (Kos island, Greece), IEEE, IEEE, 28/08/2017 2017.