

Automatic Singing Transcription

Victor Shan

December 1, 2023

Abstract

generation of new songs.

1 Introduction

Songs play an important part of society because they are conduits of cultural expression, foster emotional connection and is a great source of entertainment. Automatic Singing Transcription (AST) is a branch of Automatic Speech Recognition (ASR) can facilitate a deeper way to interact with it. In ASR, the goal is to transcribe normal speech to text, whether that’s a conversation between two people, a lecture, or a podcast. Examples of ASR include Youtube’s auto-generated captions, assistants like Siri and Alexa, and even live captioning of meetings through applications like Zoom. AST is a similar task, but instead of transcribing normal speech, AST transcribes singing. Applications of AST range from karaoke applications to Music Information Retrieval (MIR) systems. These MIR systems can be used to search for songs based on lyrics, categorize songs based on their lyrics, or even aid in the

While AST is similar to ASR, there are a few key differences that make it a more challenging task. First, singing is a more complex signal than normal speech. Singing has a more dynamic signal, with more variation in pitch, duration and vibrato [1]. For example, vowels sounds are often held for longer in singing than in normal speech like dragging the last part of “bye” in word “good-bye”. Songs also often include music in the background, which can make it more difficult to distinguish between the singer’s voice and the background music. Finally, there is a lack of large, publicly available datasets for AST [2]. This is in contrast to ASR, where there are many large datasets available, such as LibriSpeech [3]. This lack of data makes it especially difficult to train AST models.

There are many possible outputs for AST but one of the most detailed and useful outputs is the time-aligned phoneme sequence. This is a sequence of phonemes, or the smallest unit of sound in a language, that is aligned with the audio. Time-aligned phoneme sequences are useful because it can contains

enough information to reconstruct the not only the lyrics, but also how they match with the music. This output can be used to generate a karaoke application that highlights the lyrics as they are sung, or even train a model to sing covers of songs. Word-level alignment is an easier alternative but it misses the crucial information of the duration of phonemes within words. The ideal output is one that includes even more information such as the notes to sing at so that singers would be able to read and sing the song without any additional information. In this paper, we will mainly focus on the phoneme-level alignment but use word-level alignment to compare the performance of our model to other models.

2 Datasets

2.1 Requirements

The requirements of a good dataset are audio files, lyrics and timestamps. The audio files can be in any format such as mp3 or wav. The lyrics can be in any format as well, but the timestamps must be in a format that can be used to align the lyrics with the audio. The timestamps can be in the form of a phoneme sequence or word sequence but should also include the onset and offset times (beginning and end). Sentence level timestamps are too coarse and would need to be broken down into word or phoneme level timestamps. Most lyrics used by popular music services such as Spotify only use sentence/line level timestamps. Ideally, the dataset would

also be large enough to train deep neural networks and publicly available so that other researchers can use it and compare their models to existing models.

2.2 Challenges

With these requirements in mind, it is easy to see why there are so few datasets available. First, it is difficult to create a dataset with lyrics and timestamps. The lyrics must be manually transcribed and the timestamps must be manually aligned with the audio. This is a time consuming process that requires a lot of effort. Second, it is difficult to obtain the rights to use the audio files. Most songs are owned by record labels and it is difficult to obtain the rights to use these files. These two challenges make it difficult to create new datasets and is the reason why there are so few datasets available and why most of the existing datasets are small.

2.3 Dataset Augmentation

2.3.1 SpecAugment

Due to the lack of datasets, it is more important to make the most of the existing datasets. One way to do this is to augment the existing datasets. SpecAugment is a series of techniques that augment the audio spectrogram to improve the performance of ASR models. A frequent intermediate representation are Mel-Frequency Cepstral Coefficients (MFCC). These are image representations of energy at certain frequencies on a

scale that more closely matches human hearing [4]. These techniques include time warping, frequency masking and time masking [5]. Time warping is a technique that stretches or compresses the audio spectrogram in the time dimension. Frequency masking is a technique that masks a random number of frequency channels in the audio spectrogram. Time masking is a technique that masks a random number of time steps in the audio spectrogram. These techniques are used to augment the audio spectrogram before it is fed into the ASR model. [5]

2.3.2 Transforming Existing Datasets Into Pseudo-Singing Datasets

Another way to augment the existing datasets is to transform the existing datasets. This can be done by shifting the pitch, duration or vibrato [6]. The advantage of this technique is that it can also be applied to speech datasets and transform them into pseudo-singing datasets. The disadvantage is that the results will contain artifacts from the transformation based on the techniques applied. Neural network models showed an almost 15% improvement on the transformed TIMIT dataset than ones trained on the original TIMIT dataset [6].

2.3.3 Transforming Utterance level datasets into Phoneme level datasets

A technique from ASR that can also be applied to AST is transforming utter-

ance/sentence level datasets into phoneme level datasets. There are many ASR datasets that contain single utterances such as LibriSpeech and MUSDB18 [3], [7]. These datasets can be effectively transformed into phoneme level datasets by using a phoneme dictionary such as the CMU Pronouncing Dictionary [8]. This dictionary contains a mapping from words to phonemes. Using this dictionary, the utterances can be transformed into phoneme sequences. These sequences can then be used with a Connectionist Temporal Classification (CTC) loss function to train AST models [9]. CTC allows for the model to output a sequence of phonemes per time step and the repetitions are collapsed into a single phoneme as shown in Figure 1. This technique was used to allow a model to incorporate the utterance level LibriSpeech dataset into the training for a phoneme ASR model [10]. Timestamps can be retrieved from the pre-CTC output that had time-aligned phoneme classifications. The same process can be applied to AST models to allow them to use the utterance level song datasets.

2.3.4 Teacher-Student Approach

The teacher-student approach is inspired by the technique of the same name that was intended to reduce the size of large deep neural networks. The idea was first train a large deep neural network, the teacher, and then train a smaller deep neural network, the student, to mimic the teacher [11]. However, this technique can also be applied in cases of

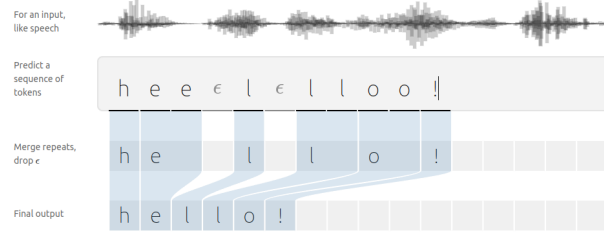


Figure 1: CTC Function collapsing repetitions [9]

low labeled data availability but high unlabeled data availability. To start off, a model would be trained on a small dataset of labeled data. Then, the model would be used to label a large dataset of unlabeled data. Finally, a new model would be trained on the newly labeled dataset [2]. This has proven to be effective for transcribing drums in music with the student model outperforming the teacher model [12].

2.4 Existing Datasets

2.4.1 TIMIT (ASR)

TIMIT is a dataset of speech recordings of 630 speakers of eight major dialects of American English with time-aligned phoneme sequences [13]. It is a popular dataset for ASR and has been used to train many ASR models. However, since this is not a singing dataset, it does not contain any of the characteristics of singing such as pitch, duration or vibrato. This is an excellent dataset to apply transformation into a pseudo-singing dataset mentioned in 2.3.2. The popularity of this dataset makes it a good candidate for applying transformations to create a pseudo-singing dataset and also for a general bench-

mark to compare against other models.

2.4.2 LibriSpeech (ASR)

LibriSpeech is a 1000 hour dataset of audio-book recordings where each recording has a matching sentence [3]. This dataset is also a popular dataset for ASR and has been used to train many ASR models including recent breakthroughs like wav2vec 2.0 [10]. This dataset is also a good candidate for training AST models because it is large and publicly available. This dataset is special because of the clarity of the audio and previous success by other models such as the wav2vec 2.0 model in detecting phonemes in this dataset when fine tuned with the TIMIT dataset [10]. This is done through the CTC technique mentioned in 2.3.3. Since speech and singing both use the same phonemes, because they are in the same language, this large dataset can be used to train base models before fine-tuning them AST models.

2.4.3 Jamendo Dataset

This dataset one of the most popular datasets for AST and has been used by many state-of-the-art AST models. It contains 20 En-

glish songs and 60 songs in other languages with word-aligned timestamped sequences [14]. This dataset is a good candidate for training AST models because it is publicly available and it’s popularity makes it an excellent benchmark to compare against other models. However, it is still a relatively small dataset and does not contain any phoneme sequences. The authors of this dataset were able to achieve a 77.8% Word Error Rate (WER) which still leaves a lot of room for improvement. **Remember to record the performance of the state of the art models on this dataset**

2.4.4 MUSDB18

MUSEDDB18 is a dataset of 150 songs with isolated vocals and accompaniment tracks [7]. This dataset is a good candidate for training AST models because it is publicly available and it has clean isolated vocals. The downside of this dataset is that it doesn’t contain any word level timestamps. However, with the CTC technique mentioned in 2.3.3, this dataset can still be used to train AST models.

2.4.5 NUS Dataset

This dataset is one of the few datasets that contains phoneme level timestamps [15]. There are 169 minutes of 20 unique English songs by 12 different people. The CMUDict was used for the phoneme vocabulary and timestamps were manually annotated. This dataset is the ideal dataset type for training

AST models due to this level of detail. It also includes a mix of slow to fast melodies and a balanced gender distribution [15].

2.4.6 Free Music Archive

The Free Music Archive (FMA) dataset is a dataset of 106,574 tracks with 161 genres [16]. This dataset is not a good candidate for directly training AST models because it does not contain any lyrics and some music may not even contain any singing at all. However, it is a good source of unlabeled songs that could be labeled through the teacher-student technique in 2.3.4. It can also provide a good source of general singing audio that can be used in the training of wav2vec2.0 models that will be described in section 3.2.

2.4.7 VocalSet

VocalSet is a 10 hour dataset of a capella singing from 20 professional singers demonstrating a variety of singing techniques [17]. This dataset is a good candidate for training Voice Activity Detection (VAD) models because it contains onset and offset timestamps for each vocal segment. This is also a good a good dataset to help train AST models to know what singing sounds like.

2.4.8 Other Datasets

Many datasets were considered but left out due to the lack of availability. Some of the most popular datasets such as Mauch’s Dataset [18] and Hansen’s Dataset [19] are not publicly available anymore. Newer

datasets such as DALI [2] and DAMP! [20] are hidden behind institutional logins and require manually requesting access.

While both the Mauch’s Dataset and Hansen’s Datasets are quite small (Mauch has 20 songs, Hansen has 9 [18]), the DALI and DAMP! datasets are much larger. The DALI dataset in particular used a version of the Teacher Student technique mentioned in 2.3.4 to label 105 songs with timestamps for the word and phoneme level [2]. The DAMP! dataset is even larger with 300x30x2 song dataset. Both of these datasets would be excellent candidates for training AST models from their size alone.

3 Related Works

3.1 HMM Based Acoustic Models

The traditional approach to AST is can be separated into a pipeline of distinct steps. The first step is to extract the features from the audio. This usually includes some form of spectral analysis such as MFCCs. The second step is to use an acoustic model to classify the features into phonemes and generate a sequence of phonemes using a Hidden Markov Model (HMM). This is the traditional approach to ASR and is also used in AST [19].

TODO: Add more detail about HMMs
 TODO: Write about Dynamic Time Warping
 TODO: Write about some results of HMMs

3.2 wav2vec 2.0 and Transfer Learning

3.2.1 wav2vec 2.0

One of the most recent breakthroughs in ASR is wav2vec 2.0 [10]. This model uses a self-supervised learning approach was used during the training of the model. Unlabeled audio was fed into the model to learn discrete speech units. These discrete speech units required the Gumbel-Softmax [21] to allow for backpropagation. The model was then fine-tuned with a linear layer and CTC loss on labeled data to perform ASR [10].

This approach was able to achieve state-of-the-art results on the LibriSpeech dataset [3] and the TIMIT dataset [13]. However, the thing that makes this approach the most promising is the fact that after pre-training on a large amount of unlabeled data, the model can be fine-tuned on a 10 minute subset of labeled data to achieve 5.2 WER on the LibriSpeech clean dataset [10]. This is very important for AST because there are so few labeled datasets available.

3.2.2 wav2vec 2.0 Transfer Learning

Using transfer learning for AST using wav2vec 2.0 was attempted in 2022 and was used to achieve state-of-the-art results on the Jamendo dataset as well as on the DALI, Hansen, Mauch, and DAMP! datasets [22]. The approach they performed did not change the first part of the model that generated the discrete speech units. Instead, they changed the last part of the model to have another

branch that outputs the probability of the current word given previous words and context. This approach achieved a 33.13 WER on the Jamendo dataset [22].

Apart from their results, there were also a few other interesting things about their approach. Since the initial wav2vec 2.0 model was trained on the LibriSpeech dataset, they wanted to make the input audio of their model as similar to the LibriSpeech dataset as possible. To do this, they used Demucs v3 [23] to separate the vocals from the accompaniment. They did not further remove any noise or singing specific features from the audio. For their labeled datasets, they used utterance level labels excluding the instrumental parts. For the output, they used character level labels.

3.3 Whisper Word-Level Alignment

4 Method

1. Preprocess datasets
2. Augment datasets
3. Create frankenstein dataset
4. Train model
5. Evaluate model
6. Label unlabeled datasets
7. train student model on newly labeled datasets
8. evaluate student model on original, manually labeled datasets
9. repeat

5 Results

5.1 WER

5.2 PER

6 Discussion

7 Future Work

7.1 Adversarial Training

7.2 Transform Singing to Speech Data

remove pitch from singing

8 Conclusion

References

- [1] C. Gupta, R. Tong, H. Li, and Y. Wang, *Semi-supervised lyrics and solo-singing alignment*, Jul. 2018.
- [2] G. Meseguer-Brocal, A. Cohen-Hadria, and G. Peeters, “Dali: A large dataset of synchronized audio, lyrics and notes, automatically created using teacher-student machine learning paradigm,” in *19th International Society for Music Information Retrieval Conference, IS-MIR*, Ed., Sep. 2018.

- [3] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210. DOI: 10.1109/ICASSP.2015.7178964.
- [4] M. Raissi. “Mel-spectrogram and mfccs.” (2021), [Online]. Available: https://www.youtube.com/watch?v=hF72sY70_IQ.
- [5] D. S. Park, W. Chan, Y. Zhang, *et al.*, “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *Interspeech 2019*, ISCA, 2019. DOI: 10.21437/interspeech.2019-2680. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2019-2680>.
- [6] A. Kruspe, “Training phoneme models for singing with ”songified” speech data,” Jan. 2015.
- [7] Z. Rafii, A. Liutkus, F.-R. Stöter, S. I. Mimilakis, and R. Bittner, *The MUSDB18 corpus for music separation*, 2017. DOI: 10.5281/zenodo.1117372. [Online]. Available: <https://doi.org/10.5281/zenodo.1117372>.
- [8] CMU, *Cmu pronouncing dictionary*, 2023. [Online]. Available: <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>.
- [9] A. Hannun, “Sequence modeling with ctc,” *Distill*, 2017, <https://distill.pub/2017/ctc>. DOI: 10.23915/distill.00008.
- [10] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, *Wav2vec 2.0: A framework for self-supervised learning of speech representations*, 2020. arXiv: 2006.11477 [cs.CL].
- [11] S. Abbasi, M. Hajabdollahi, N. Karimi, and S. Samavi, *Modeling teacher-student techniques in deep neural networks for knowledge distillation*, 2019. arXiv: 1912.13179 [cs.CV].
- [12] C.-W. Wu and A. Lerch, “Automatic drum transcription using the student-teacher learning paradigm with unlabeled music data,” Oct. 2017.
- [13] J. S. Garofolo, L. F. Lamel, W. M. Fisher, *et al.*, *Semi-supervised lyrics and solo-singing alignment*, 1993.
- [14] D. Stoller, S. Durand, and S. Ewert, *End-to-end lyrics alignment for polyphonic music using an audio-to-character recognition model*, 2019. arXiv: 1902.06797 [cs.SD].
- [15] Z. Duan, H. Fang, B. Li, K. C. Sim, and Y. Wang, “The nus sung and spoken lyrics corpus: A quantitative comparison of singing and speech,” in *2013 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*, 2013, pp. 1–9. DOI: 10.1109/APSIPA.2013.6694316.

- [16] M. Defferrard, K. Benzi, P. Vandergheynst, and X. Bresson, “FMA: A dataset for music analysis,” in *18th International Society for Music Information Retrieval Conference (ISMIR)*, 2017. arXiv: 1612 . 01840. [Online]. Available: <https://arxiv.org/abs/1612.01840>.
- [17] J. Wilkins, P. Seetharaman, A. Wahl, and B. Pardo, “Vocalset: A singing voice dataset,” in *International Society for Music Information Retrieval Conference*, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:53875542>.
- [18] MIREX, *2021:automatic lyrics transcription*, 2021. [Online]. Available: https://music-ir.org/mirex/wiki/2021:Automatic_Lyrics_Transcription.
- [19] J. K. Hansen, “Recognition of phonemes in a-cappella recordings using temporal patterns and mel frequency cepstral coefficients,” 2012. [Online]. Available: <https://api.semanticscholar.org/CorpusID:37419482>.
- [20] I. Smule, *Damp-mvp: Digital archive of mobile performances - smule multilingual vocal performance 300x30x2*, 2018. [Online]. Available: <https://zenodo.org/records/2747436>.
- [21] E. Jang, S. Gu, and B. Poole, *Categorical reparameterization with gumbel-softmax*, 2017. arXiv: 1611 . 01144 [stat.ML].
- [22] L. Ou, X. Gu, and Y. Wang, *Transfer learning of wav2vec 2.0 for automatic lyric transcription*, 2022. arXiv: 2207 . 09747 [eess.AS].
- [23] S. Rouard, F. Massa, and A. Défossez, “Hybrid transformers for music source separation,” in *ICASSP 23*, 2023.