

# SKINAUGMENT: AUTO-ENCODING SPEAKER CONVERSIONS FOR AUTOMATIC SPEECH TRANSLATION

Arya D. McCarthy<sup>\*†</sup>    Liezl Puzon<sup>†</sup>    Juan Pino<sup>†</sup>

<sup>\*</sup> Center for Language and Speech Processing, Johns Hopkins University  
<sup>†</sup> Facebook

## ABSTRACT

We propose autoencoding speaker conversion for training data augmentation in automatic speech translation. This technique directly transforms an audio sequence, resulting in audio synthesized to resemble another speaker’s voice. Our method compares favorably to SpecAugment on English–French and English–Romanian automatic speech translation (AST) tasks as well as on a low-resource English automatic speech recognition (ASR) task. Further, in ablations, we show the benefits of both quantity and diversity in augmented data. Finally, we show that we can combine our approach with augmentation by machine-translated transcripts to obtain a competitive end-to-end AST model that outperforms a very strong cascade model on an English–French AST task. Our method is sufficiently general that it can be applied to other speech generation and analysis tasks.

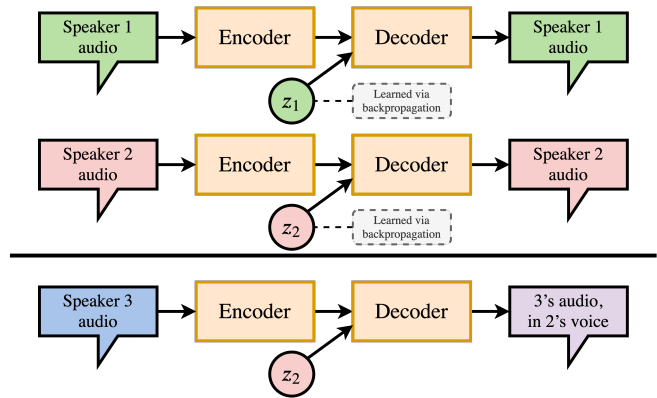
**Index Terms**— automatic speech translation, end-to-end speech translation, data augmentation, speaker normalization

## 1. INTRODUCTION

The rarity of organic training examples presents a dilemma for automatic speech translation (AST); present-day AST is a low-resource task. While **end-to-end** models seem preferable from the perspective of inference latency or error propagation, they are difficult to train to competitive levels of performance. By contrast, **cascade** models [1, 2] are not bound to audio samples and their translations. They can leverage large-scale automatic speech recognition (ASR) and machine translation (MT) training datasets. Is it possible to create an AST model with high performance while keeping the benefits of end-to-end systems?

Data augmentation is a common solution for low resource scenarios and has been explored for both ASR and AST. One of the most recent and successful data augmentation methods, SpecAugment [3], modifies the spectrogram with time warping, frequency masking and time masking. AST methods to leverage ASR and MT data include pretraining [4], multitask learning [5] and weakly supervised data augmentation [6, 7].

In this work, we generate additional audio samples without requiring transcripts, using a recent neural voice conversion



**Fig. 1.** Conditional autoencoding of two speakers’ audio (top). The latent speaker representation  $z_i$  can transform (“skin”) new audio from unseen speakers (bottom).

technique, “text-to-speech skins” [8]. Operating on the raw wav audio, it isolates the essential from the contextual aspects of speech, transferring essential aspects into a new voice [9]. We apply this method to samples of AST and ASR training data to produce new variants, in a process we call SKINAUGMENT. We additionally investigate neural speaker normalization based on the same conversion model.

We assess our proposals on English–French and English–Romanian AST tasks as well as on a low-resource English ASR task. We compare SKINAUGMENT to SpecAugment.

We find that SKINAUGMENT effectively improves the performance of end-to-end AST models, without requiring additional annotated AST or MT data. Particularly, we see BLEU gains of 2.2 on En–Fr and 3.3 on En–Ro. SKINAUGMENT outperforms both SpecAugment and a simple stochastic alteration we propose that improves SpecAugment on two AST tasks and one low-resource ASR task. However, we find no significant benefit to using SKINAUGMENT for test set normalization.

Further, we are able to produce a competitive end-to-end AST system by combining SKINAUGMENT and weak supervision from machine-translated ASR samples. This system outperforms a very competitive cascade [7] by 1.1 BLEU.

## 2. SKINAUGMENT: AUGMENTATION WITH VOICE CONVERSION BY CONDITIONED AUTOENCODING

Our speaker conversion technique [8] employs a convolutional wav-to-wav network, summarized in Figure 1. The end-to-end encoder–decoder architecture optimizes an autoencoding loss, reproducing (a shifted version of) the original input while conditioned on a latent speaker representation. This representation is learned by backpropagation while optimizing the cross-entropy  $\ell$  over  $N$  training samples  $\mathbf{x}^{(i)}$ :

$$\mathcal{L}(\theta) = \sum_{i=1}^N \ell\left(\text{dec}\left(\text{enc}(\mathbf{x}^{(i)}), \mathbf{f}_0(\mathbf{x}^{(i)}), \mathbf{z}_{s(i)}\right), \mathbf{x}^{(i)}\right). \quad (1)$$

Here,  $s$  is a function that maps training indices to speaker IDs.  $\mathbf{Z} \in \mathbb{R}^{|S| \times d}$  is a matrix with  $d$ -dimensional latent representation for each of the  $|S|$  speakers seen during training. Extracting the fundamental frequency series with  $\mathbf{f}_0$  helps to preserve the original audio’s prosody. We convert to new speakers by priming the decoder with the intended speaker’s embedding. Training does not require parallel audio recordings between speakers, nor does it require transcripts of the audio. New speakers can easily be introduced by fine-tuning the model, conditioned on a new representation.

The method achieves competitive performance on the Voice Conversion Challenge 2018 benchmark [10], despite using fewer parameters than winning systems. While the method’s value to the voice conversion task has been demonstrated [8], we show its utility for achieving superior *downstream* performance. (A related spectrogram-to-spectrogram voice converter has been applied to speech separation, trained on approximately 150 times as many hours of data [11]. In principle, this method or others could also be employed as the voice conversion subcomponent in SKINAUGMENT.)

**Augmentation Policy.** One may ask whether sheer quantity of data or its diversity contributes more to performance: Is it helpful to hear diverse variants of the same audio? Our augmentation procedure, SKINAUGMENT, lets us address this: We sample a fraction of the training data, then skin this subset into any of  $K$  arbitrarily chosen voices. In this work, we experiment with up to 16 skinned variants of the training data, sampling between 10% and 100% of the data to be skinned.

## 3. EXPERIMENTAL SETUP

### 3.1. Datasets and Evaluation

An AST dataset pairs source-language audio with a target-language translation. We experiment on two standard AST datasets: AST LibriSpeech [12] (English–French; we use the same setup as [13]) and MuST-C (English–Romanian; 432 hours) [14]. We also use AST LibriSpeech for low resource ASR. (Note that the AST LibriSpeech test sets do not correspond to the original LibriSpeech test sets.) In all cases, we use the 200+ voices in LibriSpeech to train the conversion

model; the LibriSpeech tasks let us evaluate conversions on *in-domain* audio, while the MuST-C task evaluates conversions on *out-of-domain* audio. As the AST LibriSpeech corpus’s test set is a subset of LibriSpeech’s training set, we remove all the AST LibriSpeech test set voices from LibriSpeech’s training set before training the converter.

In later experiments, we further augment the training data by translating LibriSpeech’s transcripts (removing test set occurrences [7]) with an MT system. The MT system is trained on two standard datasets: WMT16 for En–Ro (600k sentence pairs) and WMT14 for En–Fr (29 million sentence pairs).

Our En–Fr AST cascade baseline’s MT subsystem is trained on the same WMT corpora. The ASR subsystem is trained on the full LibriSpeech corpus.

For AST, we report BLEU [15] on tokenized output. (On the ASR task, the transcript is already tokenized; on the AST tasks, we tokenize translations with Moses [16].) For ASR, we use word error rate (WER), also on tokenized output.

### 3.2. Model Architecture

All of our experiments use the same mixed convolutional–recurrent end-to-end model architecture for conditional sequence generation, our focus being data augmentation techniques. (Recent work suggests that AST performance with Transformer is similar to AST performance with this style of model [17].) We use a speech encoder consisting of two non-linear layers followed by two convolutional layers and three bidirectional LSTM layers, along with a custom LSTM decoder [13, 7]. The encoder uses 40 log-scaled mel spectrogram features. We use 3 decoder layers as in [7], who report the number of parameters in each model.

SKINAUGMENT couples an off-the-shelf, fixed time-delay neural network (TDNN) encoder with a learned WaveNet decoder. Hyperparameters are as in [8].

### 3.3. Baselines

**Cascade.** We compare our data-augmented end-to-end model to a baseline cascade model. The ASR model is described in subsection 3.2, while the MT model uses a Transformer, trained on the WMT14 En–Fr parallel data. It achieves top performance on the AST LibriSpeech dataset of 21.3 BLEU [7]. We also compare to [13]’s cascade which lacks additional data.

**SpecAugment.** We also compare our end-to-end model with another popular data augmentation strategy, spectral augmentation (SpecAugment). SpecAugment adds perturbations at the feature level, whereas SKINAUGMENT operates at the raw wave level. We use the *LibriSpeech double* setting [3].

**SpecAugment- $p$ .** Further, we introduce a simple but effective variant: SpecAugment- $p$ , which applies SpecAugment to each batch with probability  $p$ . (The standard SpecAugment would thus use  $p = 1$ .) We found that SpecAugment with  $p = 0.5$  was effective in our setup.

### 3.4. Augmentation and Normalization Settings

We perform conditional generation of new data with either 8 or 16 voice conversions, applying them to 10%, 25%, 50%, and 100% of the training corpus. This creates transformed variants of our dataset in distinct (arbitrarily selected) voices. Generation is performed offline—thus not a prerequisite for inference—and is agnostic toward the AST model. While future work can explore a greater number of voices, we found this prohibitive in terms of training time.

We compare these settings to standard SpecAugment, as well as to SpecAugment- $p$  with  $p = 0.5$ , which we found to be effective.

Perhaps rather than making the AST model robust to the niceties of individual speakers’ voices, we ought to eliminate those niceties. To test the effectiveness of translation on a consistent voice, we convert the test set to entirely be of the voice of one speaker and evaluate the BLEU score separately on these single-speaker skinned test sets. We select 8 voices arbitrarily. We then produce 8 such skinned test sets with SKINAUGMENT, reporting average performance and standard deviation across the variants.

### 3.5. Machine-Translated Augmentation

Existing AST samples are rare, leading research to explore avenues for weak supervision. Among these, machine-translated transcripts of large ASR corpora dramatically increase the performance of AST models [6, 7]. We therefore translate LibriSpeech transcripts with our Transformer, then concatenate these synthetic training instances to the AST data. We apply 16 skins to 25% of the AST training data, as we found this to perform best.

### 3.6. Training Settings

We use the Adam optimizer [18] with a learning rate of 0.001 and gradient clipping of 5. The minibatch size is 96,000 frames. All experiments are conducted on 8 NVIDIA Tesla V100 GPUs. In order to compensate for the imbalance between synthetic skinned data and original data, all models are fine-tuned on the original data for 40 epochs after convergence on the augmented data. We found a consistent improvement from fine-tuning.

We decode with a beam size of 20. To balance between the data sparsity of a word-level model and the training time of a character-level model, we use a SentencePiece [19] unigram model with vocabulary size 10,000.

## 4. EXPERIMENTAL RESULTS

Results are presented in Figure 2 (En–Ro AST), Figure 3 (En–Fr AST), and Figure 4 (ASR). On all three tasks, we find that our augmentation strategy outperforms SpecAugment and SpecAugment- $p$ . We also found that SpecAugment- $p$  outperforms SpecAugment on all tasks except ASR.

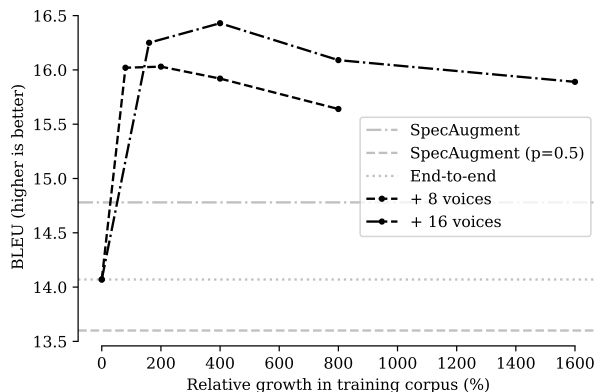


Fig. 2. English–Romanian AST with out-of-domain skins. SKINAUGMENT outperforms both SpecAugment variants.

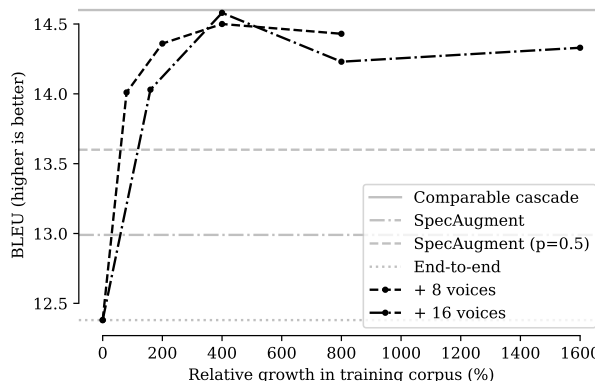


Fig. 3. English–French AST with in-domain skins. Additional synthetic data is eventually harmful.

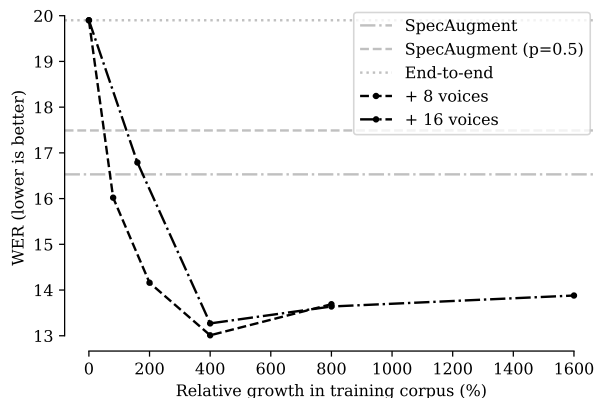


Fig. 4. SKINAUGMENT outperforms both variants of SpecAugment on the simulated low-resource ASR scenario.

SKINAUGMENT improves BLEU by 3.3 points for En–Ro and 2.2 for En–Fr over the end-to-end baseline. Our score of 14.58 matches the reported En–Fr score of [13] with a cascade model (14.6), up to their reported significant figures.

#### 4.1. Quantity versus diversity of the augmented data

How much augmented data is needed for strong performance? Does the advantage in the previous subsection come from pure quantity of data, or is the diversity of speakers advantageous? We find that the performance of the end-to-end model tracks remarkably well with the *amount* of data added, regardless of whether it comes from eight skins or sixteen, up to about twice the original size of the training data.

Beyond this, having more skins seems to be beneficial. When the amount of skinned data is  $4\times$  the size of the training set, the 16-skin model (skinning 25% of the training data) has a relative gain from the end-to-end baseline of 58%, compared to the 8-skin model (skinning 50% of the training data). Above this training set size (up to  $16\times$ ), performance begins to degrade. Nevertheless, it continues to outperform the baseline end-to-end model in our setting. For ASR and En–Ro AST, performance plateaus.

#### 4.2. Test set normalization

For all three tasks, we also skin the test set to a single training-set voice, then evaluate. The motivation is to reduce variation in test data. To avoid reporting fortuitous but unrepresentative performance from a particular voice, we consider mean BLEU and standard deviation across 8 voices. Here, we report a negative result. In every case, the score on the unmodified test set is within one standard deviation of the voice-normalized mean. In subsection 4.3, we find that translating normalized variants underperforms translating on the original audio by 0.6 BLEU on average, suggesting that the massive amount of unskinned audio obviates the benefit of skinned test data.

Our findings mesh well with [20], who found that their Cycle-GAN voice converter was harmful for test set normalization and had negligible value for data augmentation; nevertheless, in their case and ours, increased amounts of skinned data led to better performance on normalized test sets. Future work can explore whether fine-tuning to the normalization voice improves performance.

#### 4.3. Machine-Translated Data for Augmentation

Thus far, our synthetic data on the AST task has been generated by transforming original AST samples with either SKINAUGMENT or SpecAugment(-*p*). However, adding translated data as weak supervision in our low-resource scenario improves performance significantly. Table 1 shows an ablation: incorporating SKINAUGMENT, translated transcripts (“+ MT”), or both. Furthermore, we demonstrate performance when using the original AST corpus from the augmented LibriSpeech

**Table 1.** Value of machine-translated transcripts combined with SKINAUGMENT on AST LibriSpeech. We use 16 skins applied to 25% of the corpus.

Data	BLEU
AST LibriSpeech	13.24
+ SKINAUGMENT	15.22
+ MT	19.71
+ MT + SKINAUGMENT	20.19
AST LibriSpeech – AT	1.81
+ MT	21.78
+ MT + SKINAUGMENT	22.44
Cascade (with ASR and MT data)	21.31

release [12], i.e. removing the off-the-shelf automatic translations added in [13]’s dataset (“– AT”). We speculate that the abysmal performance of the baseline AST LibriSpeech – AT is due to data scarcity, and that removing automatic translations for + MT helps because they are of lower quality.

## 5. CONCLUSION

We have evaluated speaker conversions using conditioned autoencoding for AST and ASR data augmentation. A wav-to-wav CNN architecture learns latent speaker representations. Swapping in a new speaker representation converts the voice in the audio. This yields more source audio for a given example. The method is applicable to both data augmentation during training and speaker normalization for generation.

While this method relies on additional audio data to train the speaker conversion, it does not rely on transcribed text, which makes it appealing for scaling to different languages and in low-resource scenarios where annotation can be costly. SKINAUGMENT compares favorably to SpecAugment, a popular data augmentation method that operates at the feature level. We were also able to effectively combine speaker conversion data with MT-augmented ASR data. Still, when instead applied to the test set at inference time as voice normalization, we observe no significant change in BLEU.

Creating AST data by text-to-speech (TTS) synthesis of parallel text corpora has shown mixed results; while [6] found that adding a TTS system’s outputs improved performance, [7] were unable to find additional gains. The promise of SKINAUGMENT to produce variants of a given audio without transcripts suggests that it could apply to such TTS data. Future work will explore the application of this augmentation approach to improving the effectiveness of TTS data for AST.

## 6. ACKNOWLEDGMENTS

We thank Hyunbin Park, Adam Polyak, Xiaohui Zhang, and Weiyi Zheng for assistance in generating voice-converted data.

## 7. REFERENCES

- [1] Matt Post, Gaurav Kumar, Adam Lopez, Damianos Karakos, Chris Callison-Burch, and Sanjeev Khudanpur, “Improved speech-to-text translation with the Fisher and Callhome Spanish–English speech translation corpus,” in *Proc. IWSLT*, 2013.
- [2] G. Kumar, M. Post, D. Povey, and S. Khudanpur, “Some insights from translating conversational telephone speech,” in *ICASSP*, May 2014, pp. 3231–3235.
- [3] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, “SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition,” in *Proc. Interspeech 2019*, 2019, pp. 2613–2617.
- [4] Sameer Bansal, Herman Kamper, Karen Livescu, Adam Lopez, and Sharon Goldwater, “Pre-training on high-resource speech recognition improves low-resource speech-to-text translation,” in *Proc. NAACL-HLT*, Minneapolis, Minnesota, June 2019, pp. 58–68, ACL.
- [5] Ron J. Weiss, Jan Chorowski, Navdeep Jaitly, Yonghui Wu, and Zhifeng Chen, “Sequence-to-sequence models can directly translate foreign speech,” in *Proc. Interspeech 2017*, 2017, pp. 2625–2629.
- [6] Ye Jia, Melvin Johnson, Wolfgang Macherey, Ron J Weiss, Yuan Cao, Chung-Cheng Chiu, Naveen Ari, Stella Laurenzo, and Yonghui Wu, “Leveraging weakly supervised data to improve end-to-end speech-to-text translation,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 7180–7184.
- [7] Juan Pino, Liezl Puzon, Jiatao Gu, Xutai Ma, Arya D. McCarthy, and Deepak Gopinath, “Harnessing indirect training data for end-to-end automatic speech translation: Tricks of the trade,” in *Proc. IWSLT*, 2019.
- [8] Adam Polyak, Lior Wolf, and Yaniv Taigman, “TTS skins: Speaker conversion via ASR,” *CoRR*, vol. abs/1904.08983v1, 2019.
- [9] E Moulines and Y Sagisaka, “Voice conversion: state of the art and perspectives (special issue of speech communication),” 1995.
- [10] Jaime Lorenzo-Trueba, Junichi Yamagishi, Tomoki Toda, Daisuke Saito, Fernando Villavicencio, Tomi Kinnunen, and Zhenhua Ling, “The Voice Conversion Challenge 2018: Promoting development of parallel and nonparallel methods,” in *Proc. Odyssey 2018 The Speaker and Language Recognition Workshop*, 2018, pp. 195–202.
- [11] F. Biadsy, R. J. Weiss, P. J. Moreno, D. Kanvesky, and Y. Jia, “Parrotron: An end-to-end speech-to-speech conversion model and its applications to hearing-impaired speech and speech separation,” in *Proc. Interspeech*, Graz, Austria, Sept. 2019.
- [12] Ali Can Kocabiyikoglu, Laurent Besacier, and Olivier Kraif, “Augmenting Librispeech with French translations: A multimodal corpus for direct speech translation evaluation,” in *Proc. LREC*, Miyazaki, Japan, May 2018, ELRA.
- [13] Alexandre Bérard, Laurent Besacier, Ali Can Kocabiyikoglu, and Olivier Pietquin, “End-to-end automatic speech translation of audiobooks,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6224–6228.
- [14] Mattia Antonino Di Gangi, Roldano Cattoni, Luisa Benvivoglio, Matteo Negri, and Marco Turchi, “MuST-C: a multilingual speech translation corpus,” in *Proc. NAACL-HLT*, Minneapolis, MN, USA, June 2019, ACL.
- [15] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu, “BLEU: a method for automatic evaluation of machine translation,” in *Proc. ACL*, Philadelphia, Pennsylvania, USA, July 2002, pp. 311–318, ACL.
- [16] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst, “Moses: Open source toolkit for statistical machine translation,” in *Proc. ACL*, Prague, Czech Republic, June 2007, pp. 177–180, ACL.
- [17] Shigeki Karita, Nanxin Chen, Tomoki Hayashi, Takaaki Hori, Hirofumi Inaguma, Ziyang Jiang, Masao Someki, Nelson Enrique Yalta Soplín, Ryuichi Yamamoto, Xiaofei Wang, Shinji Watanabe, Takenori Yoshimura, and Wangyou Zhang, “A comparative study on transformer vs RNN in speech applications,” *CoRR*, vol. abs/1909.06317, 2019.
- [18] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *Proc. ICLR*, 2015.
- [19] Taku Kudo and John Richardson, “SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” in *Proc. EMNLP*, Brussels, Belgium, Nov. 2018, pp. 66–71, ACL.
- [20] Gokce Keskin, Tyler Lee, Cory Stephenson, and Oguz H. Elibol, “Measuring the effectiveness of voice conversion on speaker identification and automatic speech recognition systems,” 2019.