

# ON PROSODY MODELING FOR ASR+TTS BASED VOICE CONVERSION

Wen-Chin Huang<sup>1</sup>, Tomoki Hayashi<sup>1</sup>, Xinjian Li<sup>2</sup>, Shinji Watanabe<sup>2</sup>, Tomoki Toda<sup>1</sup>

<sup>1</sup>Nagoya University, Japan

<sup>2</sup>Carnegie Mellon University, USA

## ABSTRACT

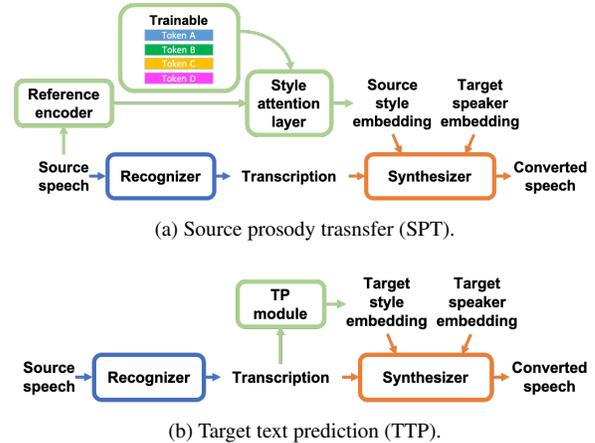
In voice conversion (VC), an approach showing promising results in the latest voice conversion challenge (VCC) 2020 is to first use an automatic speech recognition (ASR) model to transcribe the source speech into the underlying linguistic contents; these are then used as input by a text-to-speech (TTS) system to generate the converted speech. Such a paradigm, referred to as ASR+TTS, overlooks the modeling of prosody, which plays an important role in speech naturalness and conversion similarity. Although some researchers have considered transferring prosodic clues from the source speech, there arises a speaker mismatch during training and conversion. To address this issue, in this work, we propose to directly predict prosody from the linguistic representation in a target-speaker-dependent manner, referred to as target text prediction (TTP). We evaluate both methods on the VCC2020 benchmark and consider different linguistic representations. The results demonstrate the effectiveness of TTP in both objective and subjective evaluations.

**Index Terms**— voice conversion, automatic speech recognition, text-to-speech, prosody, global style token

## 1. INTRODUCTION

In voice conversion (VC), one aims to convert the speech from a source to that of a target without changing the linguistic content [1, 2]. From an information perspective, the goal of VC is to extract the spoken contents from the source speech and then synthesize the converted speech from the extracted contents with the identity of the target speaker. An ideal VC system would then consist of two components, the recognition module and the synthesis module, where the two components perform the above-mentioned respective actions. Such a paradigm can be directly realized by cascading an automatic speech recognition (ASR) model and a text-to-speech (TTS) system, which we refer to as ASR+TTS. In the latest voice conversion challenge 2020 (VCC2020) [3], ASR+TTS was adopted as one of the baseline systems [4], and the top performing system also implemented such a framework [5], showing state-of-the-art performance in terms of both naturalness and similarity.

Despite the promising results, the conversion and modeling of prosody in the ASR+TTS paradigm are often ignored. Prosody is a combination of several fundamental prior components in speech, such as pitch, stress, and breaks, and it impacts subsequent high-level characteristics including emotion and style. From the information perspective described in the previous paragraph, the synthesis module in ASR+TTS is responsible for recovering all information discarded by the recognition module. In the text-based ASR+TTS baseline system [4], since the mapping from text to prosody is one-to-many, a specific prosody or style is unpredictable from text, so the TTS model can only implicitly model the prosody pattern from



**Fig. 1:** Illustration of the conversion processes for the two prosody modeling techniques for ASR+TTS-based VC that are examined in this work.

the statistical properties of the training data, resulting in a collapsed, averaged prosodic style.

A closely related research field in which one aim is to control the variability of speech is expressive speech synthesis. One of the most widely used methods is the use of the global style token (GST) [6], with which a reference speech is encoded into a fixed dimensional embedding [7] represented as a weighted sum of a set of predefined style tokens. A TTS model equipped with GST (GST-TTS) is therefore formulated as a conditional generative model given the text, speaker identity, and prosody encoding from the reference speech. Many have generalized such a framework to variational autoencoders (VAEs) [8, 9] to increase the generalizability of the learnt embedding space. Another line of work focuses on fine-grained prosody control by learning variable-length prosody embeddings [10, 11].

The above-mentioned approach was first applied to VC in [12], and it will be referred to as *source prosody transfer (SPT)*. As illustrated in Figure 1a, the source speech was used as input by the reference encoder to generate a global prosody embedding such that the prosody of the converted speech follows the source. Such a process is also termed *same-text prosody transfer* in the expressive TTS literature [7] since the reference speech contains the same linguistic contents as the input of the TTS, which was also derived from the source speech. Several top performing teams in VCC2020 extended SPT by adopting a variational reference encoder augmented with a speaker adversarial layer to help with disentanglement [5, 13], showing the competitiveness of SPT.

Nonetheless, an ablation study in [5] showed that SPT did not bring about any significant improvements on task 1. We suspect

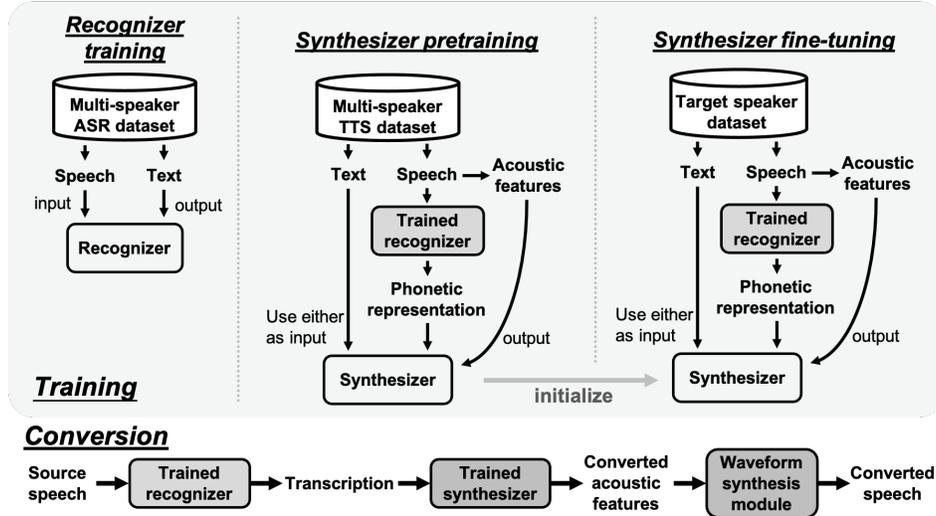


Fig. 2: Training and conversion procedures in ASR+TTS-based VC.

that SPT can be a sub-optimal strategy for prosody modeling in ASR+TTS-based VC. First, the target-speaker-dependent TTS training causes a mismatch between training and conversion, because the speech of the target speaker is used as input to the reference encoder during training but that of the source is used during conversion. A speaker adversarial classifier can alleviate this issue [5, 13] but requires careful hyperparameter tuning. Second, there are scenarios where SPT is not desired, such as emotion VC or accent conversion.

In this work, we examine two prosody modeling methods for ASR+TTS-based VC. In addition to SPT, we propose a novel technique, which we refer to as *target text prediction (TTP)*. We borrow the idea from [14] and train a text prediction (TP) module to generate the prosody embedding from the text derived from the source speech. Figure 1b illustrates this process. The TP module is first pretrained with a GST-TTS on a multispeaker dataset, and further fine-tuned in a target-speaker-dependent manner. As a result, TTP does not suffer from a mismatch between training and conversion unlike SPT. Our contributions in this work are as follows.

- We propose TTP, a new prosody modeling technique for ASR+TTS-based VC, that predicts prosody in a target-speaker-dependent manner.
- We apply two prosody modeling methods, namely, SPT and TTP, to three ASR-TTS systems that differ in representation, and evaluate them in the two tasks in VCC2020. The results show that TTP constantly outperforms SPT.

## 2. VOICE CONVERSION BASED ON ASR+TTS

### 2.1. Overall framework and conversion process

The ASR+TTS-based VC system in this work is built upon the baseline system for VCC2020 [4]. As depicted in Figure 2, the system consists of three modules: a speaker-independent recognizer, a target-speaker-dependent synthesizer, and a neural vocoder that generates the final speech waveform. Starting from the source speech  $\mathbf{X}$ , the recognizer first extracts the spoken contents:  $\hat{\mathbf{Y}} = \text{Recognizer}(\mathbf{X})$ . The synthesizer takes the transcription and

synthesizes the converted acoustic features:  $\hat{\mathbf{X}} = \text{Synthesizer}(\hat{\mathbf{Y}})$ <sup>1</sup>. The neural vocoder finally uses the converted acoustic features as input to reconstruct the waveform.

### 2.2. Intermediate representation

The spoken contents,  $\mathbf{Y}$ , can be any intermediate representation. In addition to using text as in [4], in this work we also evaluate the following two representations. Note that the goal of this work is to examine the effectiveness of the prosody modeling techniques over the three representations, but not to provide a fair comparison among them.

**Bottleneck feature (BNF):** Speaker-independent, frame-level phonetic bottleneck features derived from an ASR model provide strong clues about the content. Such features were first used in [15], where the phonetic posteriorgram (PPG), a time-versus-class matrix representing the posterior probabilities of each phonetic class (referring to a word, phone, or senone) for each frame, was used. In this work, as in [16], we use the hidden representations before the last softmax layer in an ASR model.

**VQW2V:** Features such as BNF are derived from an ASR model, which requires supervision using labels; this increases the cost of building such a system, especially in low-resource or cross-lingual settings. Alternatively, several studies [17–19] adopted self-supervised representations that do not require any label during training while still being speaker-independent and framewise. In this work, we adopt the vector-quantized wav2vec (VQW2V) [20].

### 2.3. Training

**ASR:** A multispeaker dataset  $\mathbf{D}_{\text{ASR}}$  ensures the speaker independence of the recognizer.

**TTS:** Synthesizer training involves a pretraining and a fine-tuning stage. Pretraining is performed on a multispeaker TTS dataset  $\mathbf{D}_{\text{TTS}}$ , which is followed by fine-tuning on the limited target speaker dataset  $\mathbf{D}_{\text{trg}}$ . This is a common practice in building modern neural TTS models, as pretraining ensures stable quality and fine-tuning retains high speaker similarity [21]. Such a training strategy allows for

<sup>1</sup>In this paper, the terms *Recognizer* and *Synthesizer* are used interchangeably with *ASR* and *TTS*, respectively

training on even approximately 5 minutes of data. To train a text-based synthesizer, human-labeled text is used. However, since annotating ground-truth BNF and VQW2V is impossible, we use the trained recognizer to extract the representations for synthesizer training.

The ASR and TTS models adopt sequence-to-sequence (seq2seq) structures, which were shown to improve conversion similarity by modeling the long-term dependencies in speech. Note that the two models can be separately trained and thus benefit from advanced techniques and a wide variety of datasets in their own fields.

### 3. PROSODY MODELING IN ASR+TTS-BASED VOICE CONVERSION

We examine two prosody modeling techniques in this work, namely, SPT and TTP. In this section, we first introduce GST and TP, which are the two fundamental building blocks, then describe the training processes in SPT and TTP. We omit the descriptions of the conversion procedures as have been discussed in Section 1.

#### 3.1. Global style tokens and text prediction

GST-TTS is aimed at using a *global style embedding*<sup>2</sup> (in contrast to *fine-grained* embeddings in [10]) encoded from a reference speech to capture residual attributes not specified by other input streams including text,  $\mathbf{Y}$ , and speaker label,  $s$  [6]. Formally, given a training sample  $\{(\mathbf{X}, \mathbf{Y}, s)\}$  from  $\mathbf{D}_{\text{TTS}}$ , where  $\mathbf{X}$  means a speech utterance, training involves minimizing the following L1-loss:

$$\mathcal{L}_{\text{GST}} = \|\mathbf{X} - \text{Synthesizer}(\mathbf{Y}, s, \text{RefEnc}(\mathbf{X}))\|_1, \quad (1)$$

where the RefEnc function consists of a reference encoder, a style attention layer, and a set of trainable tokens called GSTs. Specifically, the reference encoder first takes the reference speech as input and generates a fixed dimensional output vector. It is then used as the query to the attention layer to produce a set of weights over the pre-defined GSTs. The final style embedding is essentially the weighted sum of the GSTs and is fused with the hidden states of the TTS encoder.

During inference, the original GST-TTS requires either the reference speech or manual setting of a set of weights. To operate solely given the text input, the use of an extra TP module that takes the hidden states of the TTS encoder as input to approximate either the ground-truth style weights or the style embedding extracted from the target speech was proposed in [14]. The training objective of the TP module can be formed by rewriting Equation 1 as

$$\mathcal{L}_{\text{TP}} = \|\mathbf{X} - \text{Synthesizer}(\mathbf{Y}, s, \text{TP}(\mathbf{Y}))\|_1. \quad (2)$$

The TP module can be trained jointly with the GST-TTS by using a stop-gradient operator. It was shown that the TP module can generate natural speech that well reflects the variations of the training data using the text input alone [14].

#### 3.2. Source prosody transfer

SPT was first proposed in [12] and further utilized in subsequent works [5, 13, 22, 23]. The training process of SPT involves two stages.

<sup>2</sup>The terms such as *style token* and *style embedding* are terms from the original paper [6], and readers should note they are not restricted to style but include many other factors.

**Table 1:** Summary of the data conditions in VCC2020.

| Task   | Training phase     |                              | Conversion phase   |                    |
|--------|--------------------|------------------------------|--------------------|--------------------|
|        | Source             | Target                       | Source             | Converted          |
| Task 1 | 70 Eng. utterances | 70 Eng. utterances           | 25 Eng. utterances | 25 Eng. utterances |
| Task 2 |                    | 70 Man./Ger./Fin. utterances |                    |                    |

1. Pretrain GST-TTS on  $\mathbf{D}_{\text{TTS}}$  using  $\mathcal{L}_{\text{GST}}$ , as in Equation 1.
2. Finetune GST-TTS on  $\mathbf{D}_{\text{urg}}$  using  $\mathcal{L}_{\text{GST}}$ , as in Equation 1.

The conversion process of SPT can be formulated as

$$\hat{\mathbf{X}} = \text{Synthesizer}(\text{Recognizer}(\mathbf{X}), s_{\text{urg}}, \text{RefEnc}(\mathbf{X})). \quad (3)$$

Comparing Equation 1 with Equation 3, one may find that the inputs to RefEnc are different. The speech from the target speaker is used in training, whereas the source speech is used during conversion. As discussed in Section 1, if RefEnc is overly fine-tuned towards the target speaker, SPT can suffer from a speaker mismatch problem. Freezing RefEnc during fine-tuning may be a remedy, but as we will show in the experiment section, this does not alleviate the problem.

#### 3.3. Target text prediction

We propose TTP to tackle the sub-optimal problem mentioned previously. The training process of TTP involves three stages.

1. Pretrain GST-TTS on  $\mathbf{D}_{\text{TTS}}$ .
2. Pretrain TP on  $\mathbf{D}_{\text{TTS}}$  using  $\mathcal{L}_{\text{TP}}$ , as in Equation 2 with all other model parameters fixed<sup>3</sup>. Note that the GST module is no longer required from this stage.
3. Fine-tune TP and TTS on  $\mathbf{D}_{\text{urg}}$  using  $\mathcal{L}_{\text{TP}}$ , as in Equation 2.

The conversion process of TTP can be formulated as

$$\hat{\mathbf{X}} = \text{Synthesizer}(\text{Recognizer}(\mathbf{X}), s_{\text{urg}}, \text{TP}(\text{Recognizer}(\mathbf{X}))). \quad (4)$$

Comparing Equation 2 with Equation 4, we see the speaker mismatch can be avoided since the TP module uses  $\mathbf{Y}$  as input, whose speaker independence is ensured by the recognizer. TTP is therefore more robust than SPT, as we will show in the experiments.

## 4. EXPERIMENTAL SETTINGS

### 4.1. Data

We used the VCC2020 dataset [3], which contained two tasks in our evaluation. The data conditions are summarized in Table 1. Both tasks share the same two source English male and female speakers whose data were not used. There were two target male and female speakers of English in task 1 whereas in task 2 there were one male and one female speaker each of Finnish, German, and Mandarin. During conversion, the source speaker’s voice in the source language was converted as if it was uttered by the target speaker while keeping the linguistic contents unchanged. For each target speaker, 70

<sup>3</sup>This is an objective theoretically equivalent to the one in the original model [14], where the loss of the weights or embedding between those predicted from the encoder states or the speech is optimized.

**Table 2:** Objective evaluation results of task 1. Bold font indicates the best performance in the same representation. The x mark indicates no freezing of the GST-related modules, and the check mark the opposite.

| System   | Text        |              |            |             | VQW2V       |              |             |             | BNF         |        |            |             |
|----------|-------------|--------------|------------|-------------|-------------|--------------|-------------|-------------|-------------|--------|------------|-------------|
|          | MCD         | F0RMSE       | CER        | WER         | MCD         | F0RMSE       | CER         | WER         | MCD         | F0RMSE | CER        | WER         |
| Baseline | 6.49        | 29.80        | 14.4       | 24.0        | <b>7.13</b> | 31.34        | <b>10.5</b> | <b>17.8</b> | <b>6.86</b> | 30.32  | 7.9        | 13.4        |
| SPT (X)  | 6.50        | 30.89        | 7.7        | 14.6        | 7.19        | <b>31.27</b> | 15.0        | 22.5        | 7.53        | 32.17  | 7.9        | 13.5        |
| SPT (✓)  | 6.71        | 31.82        | 7.9        | 15.1        | 7.35        | 32.46        | 13.5        | 22.6        | 7.59        | 31.92  | <b>7.7</b> | <b>13.2</b> |
| TTP      | <b>6.36</b> | <b>29.72</b> | <b>7.4</b> | <b>13.1</b> | <b>7.13</b> | 31.49        | 11.6        | 18.7        | 6.90        | 29.99  | <b>7.7</b> | <b>13.2</b> |

utterances in their respective languages and contents were provided. The 25 test sentences for evaluation were shared for tasks 1 and 2.

All recognizers were trained with the LibriSpeech dataset [24]. For the multispeaker TTS dataset, we used the “clean” subsets LibriTTS dataset [25], except in the text-based system for task 2, where we merged open-source single-speaker TTS datasets in Finnish [26], German [27], and Mandarin [28], as in [4]. For each of the two tasks, a separate neural vocoder was trained with the training data of the source and target speakers.

## 4.2. Implementation

The system was implemented using ESPnet, a well-developed open-source end-to-end (E2E) speech processing toolkit [29, 30]. Following [4], the ASR model for the text-based system was based on the Transformer [31–33] with joint CTC/attention loss [34], and a RNN-based language model for decoding<sup>4</sup>. The ASR model for BNF extraction was based on TDNNF-HMM [35], where we concatenated 40-dimensional MFCCs and 400-dimensional i-vectors as input. For VQW2V, we used the publicly available pretrained model provided by fairseq [36]<sup>5</sup>, as in [17].

All synthesizers map their respective inputs to 80-dimensional mel filterbanks with 1024 FFT points and a 256-point frame shift (16 ms). The x-vector [37] was used as the speaker embedding, and we used the pretrained model provided by Kaldi<sup>6</sup>. The average of all x-vectors of the training utterances of each speaker was used during inference. Synthesizers with discrete input including text and VQW2V had a Transformer-TTS architecture [38] with detailed settings [4, 17]. For the BNF-based synthesizer, we adopted the Voice Transformer Network (VTN) [39, 40] and followed the official implementation<sup>7</sup>. For the neural vocoder, we adopted the Parallel WaveGAN (PWG) [41] and followed the open-source implementation<sup>8</sup>.

## 5. EXPERIMENTAL EVALUATIONS

### 5.1. Objective evaluation

We carried out three types of objective evaluation metric. The character/word error rates (CER/WER) from an off-the-shelf ASR sys-

<sup>4</sup>Code and pretrained models from the official implementation on ESPnet: <https://github.com/espnet/espnet/tree/master/egs/vcc20>

<sup>5</sup><https://github.com/pytorch/fairseq/tree/master/examples/wav2vec>

<sup>6</sup><https://kaldi-asr.org/models/m8>

<sup>7</sup><https://github.com/espnet/espnet/tree/master/egs/arctic/vcl>

<sup>8</sup><https://github.com/kan-bayashi/ParallelWaveGAN>

tem is not only an estimate of intelligibility but also a strong indicator of quality in VC, as shown in [42]. Our ASR engine was Transformer-based [32] and trained on LibriSpeech. 24-dimensional mel cepstrum distortion (MCD) and the F0 root mean square error (F0RMSE) were reported to assess the spectral distortion and prosody conversion accuracy, where the calculation of both metrics were based on the WORLD vocoder [43]. Note that in task 2, only CER/WER were reported, since access to the ground truth speech was not available. Model selection was based on MCD and CER in tasks 1 and 2, respectively, and the best performing models generated the samples for the subjective test.

#### 5.1.1. Results of task 1

Table 2 shows the objective results of task 1. For the text-based system, introducing prosody modeling significantly improved CER/WER. We suspect that the improved prosody pattern enabled the ASR model to easily recognize the contents correctly. Compared with SPT, TTP was shown to be more effective as it outperformed all systems. However, for frame-level features such as BNF and VQW2V, improvements from those of SPT and TTP were not significant. This result suggests that the fine-grained representations already carry abundant prosody information that is not discarded in the recognition process, such that it cannot be properly modeled by GST.

We also found that freezing the GST modules during fine-tuning degraded the performance. One possible reason is that the GSTs capture not only prosody but also other residual attributes such as channel and noise, and such mismatch must be taken care of through the fine-tuning process. Finally, compared with SPT without freezing, TTP constantly yielded a superior performance over all representations, demonstrating its superior effectiveness.

#### 5.1.2. Results of task 2

Table 3 shows the objective results of task 2. First, for the text-based system, SPT improved the performance for German and Mandarin target speakers, and TTP was beneficial for Mandarin speakers. By listening to several samples, we found that both SPT and TTP helped the model insert proper short pauses, which significantly increased intelligibility. On the other hand, both methods degraded the performance for the Finnish speakers; we assume that the improper text preprocessing for Finnish [4] made both acoustic and prosody modeling difficult. For the framewise representations, we conclude that neither SPT nor TTP much affected the intelligibility. Note that freezing GST in SPT did not cause degradation as it did in task 1.

**Table 3:** Objective evaluation results of task 2. Bold font indicates the best performance in the same representation. The x mark indicates not freezing the GST-related modules, and the check mark the opposite.

| System   | Lang. | Text        |             | VQW2V      |             | BNF        |             |
|----------|-------|-------------|-------------|------------|-------------|------------|-------------|
|          |       | CER         | WER         | CER        | WER         | CER        | WER         |
| Baseline | Fin.  | 39.1        | 67.1        | 12.6       | 21.9        | 9.4        | 16.1        |
|          | Ger.  | 11.0        | 18.8        | 9.3        | 15.5        | 7.8        | 13.6        |
|          | Man.  | 13.3        | 23.1        | 8.2        | 16.1        | 8.4        | 14.0        |
|          | Avg.  | <b>21.1</b> | <b>36.3</b> | 10.0       | 17.8        | 8.5        | 14.5        |
| SPT (✗)  | Fin.  | 47.8        | 79.9        | 12.5       | 19.9        | 8.7        | 14.6        |
|          | Ger.  | 7.9         | 14.7        | 9.8        | 17.7        | 9.2        | 15.8        |
|          | Man.  | 8.6         | 15.3        | 8.3        | 16.6        | 8.6        | 14.8        |
|          | Avg.  | 21.4        | 36.6        | 10.2       | 18.0        | 8.8        | 15.0        |
| SPT (✓)  | Fin.  | 47.4        | 80.7        | 12.3       | 21.1        | 8.9        | 15.3        |
|          | Ger.  | 8.6         | 15.9        | 10.2       | 19.0        | 8.6        | 15.4        |
|          | Man.  | 8.3         | 15.8        | 8.2        | 15.1        | 8.1        | 14.2        |
|          | Avg.  | 21.4        | 37.4        | 10.2       | 18.4        | 8.5        | 14.9        |
| TTP      | Fin.  | 51.6        | 80.5        | 10.8       | 19.9        | 8.9        | 15.0        |
|          | Ger.  | 11.3        | 20.4        | 9.9        | 17.2        | 8.0        | 13.9        |
|          | Man.  | 8.8         | 17.1        | 7.9        | 15.3        | 7.9        | 13.6        |
|          | Avg.  | 23.9        | 39.3        | <b>9.5</b> | <b>17.5</b> | <b>8.2</b> | <b>14.2</b> |

## 5.2. Subjective evaluation

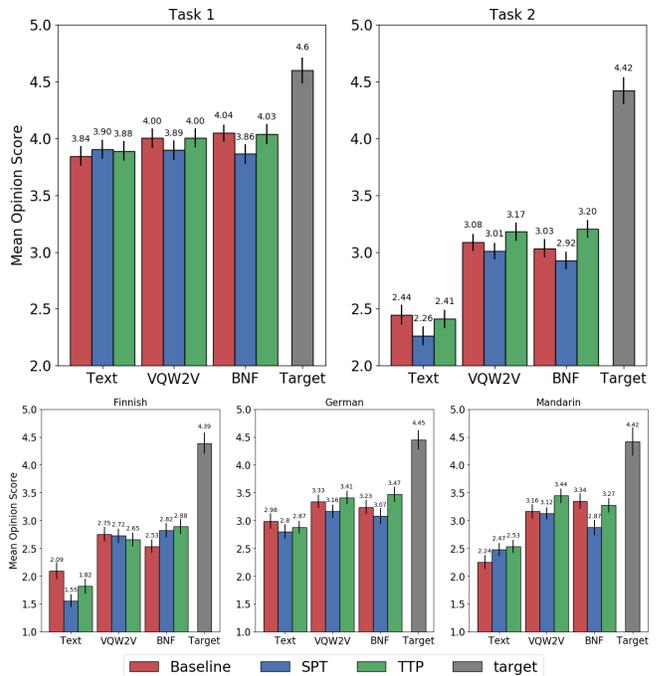
We followed the subjective evaluation methodology described in [3] and evaluated two aspects: naturalness and speaker similarity. For naturalness, participants were asked to evaluate the naturalness of the speech by the mean opinion score (MOS) test on a five-point scale. For conversion similarity, each listener was presented a natural target speech and a converted speech and asked to judge whether they were produced by the same speaker on a four-point scale.

For each system, five random utterances were chosen for each conversion pair. In the naturalness test, recordings of the target speakers were also included and served as the upper bound. In the similarity test for task 2, following [3], we selected three English recordings and two L2 language recordings as the natural reference for the five converted utterances. All subjective evaluations were performed using the open-source toolkit [44] that implements the ITU-T Recommendation P.808 [45] for subjective speech quality assessment in a crowd using the Amazon Mechanical Turk (Mturk) and screens the obtained data for unreliable ratings. We recruited more than 100 listeners from the United States and had each sample rated by five different participants on average. Note that to reduce cost, we eliminated SPT systems that freeze GST in the listening tests, since they yield inferior performance compared with systems that do not freeze GST, as shown in Section 5.1. Audio samples are available online<sup>9</sup>.

### 5.2.1. Naturalness test

Figure 3 shows the naturalness results. We focus on task 1 first. For text-based systems, contrary to the findings in Section 5.1.1, SPT

<sup>9</sup><https://unilight.github.io/Publication-Demos/publications/prosody-asr-tts-vc/index.html>



**Fig. 3:** Naturalness MOS plots. Top row: taskwise results. Bottom row: Language breakdown for results of task 2.

and TTP results were comparable to the baseline. For frame-level features, SPT largely degraded the performance, whereas TTP could compensate for the degradation but only to make it on par with the baseline. We can still conclude that TTP was superior to SPT for task 1, but no significant improvements were observed when compared with the baseline. These results are somewhat consistent with the findings in [5].

For task 2, all systems showed performance degradation brought about by SPT, and TTP significantly outperformed SPT for all representations. For systems based on frame-level features, TTP could even outperform the baseline and was also comparable to the text-based system. To investigate this gap, we plotted the breakdown per target language in the bottom graphs of Figure 3. We suspect that the relative performance change was again correlated with the text preprocessing, as stated in Section 5.1.2. As in [4], thanks to the open-source community, we utilized G2P tools to convert both English and Mandarin text into phonemes, resulting in a better acoustic model and a larger improvement brought about by TTP. On the other hand, because of the lack of linguistic knowledge, characters were used for Finnish and German, resulting in degradation when combined with TTP. We thus conclude that TTP is an effective method for improving naturalness in task 2, if the input representation is properly processed.

### 5.2.2. Similarity test

Figure 4 shows the similarity results. In task 1, both SPT and TTP demonstrated a comparable performance for the text and VQW2V-based systems, whereas SPT performed significantly worse in the BNF-based system. As for task 2, all three representations yielded a similar trend such that no obvious performance gain could be distinguished for any prosody modeling method. Although it is widely believed that prosody is highly correlated to speaker identity, here

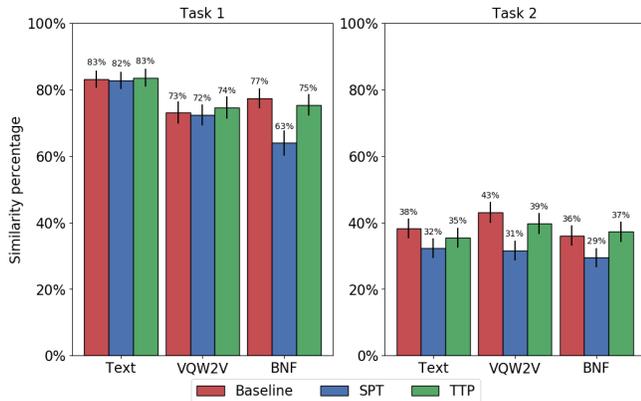


Fig. 4: Taskwise similarity results.

we could only conclude that prosody modeling only improved naturalness.

On the other hand, for situations where the L2 recordings were used as the reference speech, it was necessary to “imagine” the target speaker’s voice in English in order to judge the similarity to the converted speech. Under such a subjective evaluation protocol, we suspect that listeners could only resort to coarse prosodic clues such as the overall pitch level and timbre, which were not affected by the prosody modeling techniques.

### 5.3. Visualization of style embeddings

We visualized the style embedding space of the GST-TTS models using different representations that are pretrained on LibriTTS (which is the resulting model of stage 1 in Section 3.2). We input the training utterances of the 10 target speakers to obtain the style embeddings, and performed dimensional reduction by the t-SNE method [46] for visualization. As a reference, we also visualized the X-vectors. Finally, we colored the dots with respect to each target speaker. The resulting plots are depicted in Figure 5.

First, it can be observed that all representations (text, VQW2V, and BNF) were somehow clustered with respect to speakers. As one possible reason may be that prosody is speaker-dependent [47], we may also infer that the use of X-vectors for speaker embeddings cannot fully capture all speaker-related variations, such that the GST-TTS model relies on the help of the style embedding. This justifies the use of the speaker adversarial classifier [5, 13].

We further observed that the order of visual *degree of clusteriness* (how close together the dots of a speaker are) from low to high is text, VQW2V, BNF, and X-vector. We suggest that the degree of clusteriness reflects the amount of variation that style embeddings can explain. As we know that text contains the least prosody information, we could infer that owing to the frame-level nature, representations such as VQW2V and BNF already contained much prosody information such that there was not much left for the style embeddings to capture. Nonetheless, from the improvements of naturalness in task 2, as described in Section 5.2.1, we note that prosody modeling is still effective in capturing residual information.

## 6. CONCLUSIONS

We examined two prosody modeling methods, namely, SPT and TTP, for ASR+TTS-based VC. Whereas SPT had already been ap-

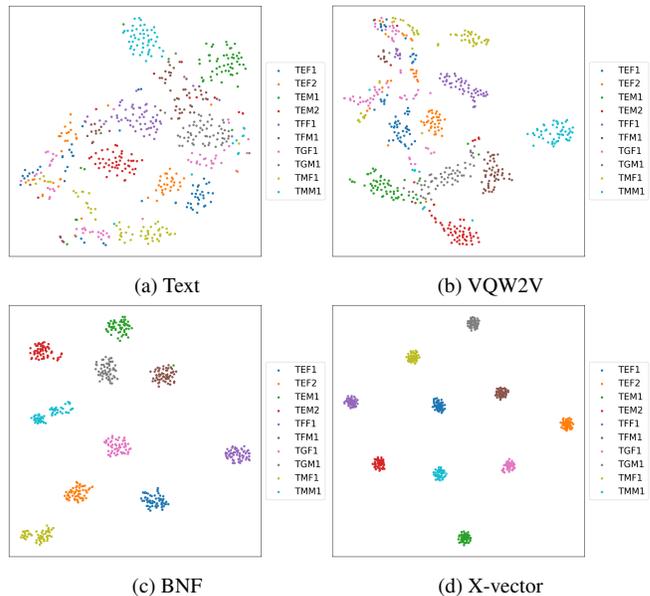


Fig. 5: Visualizations of style embeddings and X-vectors of training utterances of the 10 target speakers in VCC2020. Style embeddings were extracted using GST-TTS models pretrained on LibriTTS using different input representations. Each dot represents the embedding of one utterance.

plied by several top systems in VCC2020, TTP was newly proposed in this work, with the motivation to prevent the mismatch between training and conversion in SPT. We conducted experiments on the VCC2020 benchmark and considered three intermediate representations: text, BNF, and VQW2V. Results showed that TTP was consistently superior to SPT and could outperform the baseline with proper input representation. Finally, the visualization of the style embeddings learnt with different representations shed light on the limited improvements brought about by SPT and TTP.

## 7. ACKNOWLEDGEMENTS

This work was partly supported by JSPS KAKENHI Grant Number 21J20920 and JST CREST Grant Number JPMJCR19A3, Japan. We would also like to thank Yu-Huai Peng and Hung-Shin Lee from Academia Sinica, Taiwan, for training the BNF extractor.

## 8. REFERENCES

- [1] Y. Stylianou, O. Cappe, and E. Moulines, “Continuous probabilistic transform for voice conversion,” *IEEE TSAP*, vol. 6, no. 2, pp. 131–142, 1998.
- [2] T. Toda, A. W. Black, and K. Tokuda, “Voice Conversion Based on Maximum-Likelihood Estimation of Spectral Parameter Trajectory,” *IEEE TASLP*, vol. 15, no. 8, pp. 2222–2235, 2007.
- [3] Y. Zhao, W.-C. Huang, X. Tian, J. Yamagishi, R. K. Das, T. Kinnunen, Z. Ling, and T. Toda, “Voice Conversion Challenge 2020 - Intra-lingual semi-parallel and cross-lingual voice conversion -,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 80–98.

- [4] W.-C. Huang, T. Hayashi, S. Watanabe, and T. Toda, “The Sequence-to-Sequence Baseline for the Voice Conversion Challenge 2020: Cascading ASR and TTS,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 160–164.
- [5] J.-X. Zhang, L.-J. Liu, Y.-N. Chen, Y.-J. Hu, Y. J., Z.-H. Ling, and L.-R. Dai, “Voice Conversion by Cascading Automatic Speech Recognition and Text-to-Speech Synthesis with Prosody Transfer,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 121–125.
- [6] Y. Wang, D. Stanton, Y. Zhang, RJ-Skerry Ryan, E. Battenberg, J. Shor, Y. Xiao, Y. Jia, F. Ren, and R. A. Saurous, “Style Tokens: Unsupervised Style Modeling, Control and Transfer in End-to-End Speech Synthesis,” in *Proc. ICML*, 2018, pp. 5180–5189.
- [7] RJ Skerry-Ryan, E. Battenberg, Y. Xiao, Y. Wang, D. Stanton, J. Shor, R. Weiss, R. Clark, and R. A. Saurous, “Towards End-to-End Prosody Transfer for Expressive Speech Synthesis with Tacotron,” in *Proc. ICML*, 2018, pp. 4693–4702.
- [8] Y.-J. Zhang, S. Pan, L. He, and Z.-H. Ling, “Learning Latent Representations for Style Control and Transfer in End-to-end Speech Synthesis,” in *Proc. ICASSP*, 2019, pp. 6945–6949.
- [9] W.-N. Hsu, Y. Zhang, R. Weiss, H. Zen, Y. Wu, Y. Cao, and Y. Wang, “Hierarchical Generative Modeling for Controllable Speech Synthesis,” in *Proc. ICLR*, 2019.
- [10] Y. Lee and T. Kim, “Robust and Fine-grained Prosody Control of End-to-end Speech Synthesis,” in *Proc. ICASSP*, 2019, pp. 5911–5915.
- [11] V. Klimkov, S. Ronanki, J. Rohnke, and T. Drugman, “Fine-Grained Robust Prosody Transfer for Single-Speaker Neural Text-To-Speech,” in *Proc. Interspeech*, 2019, pp. 4440–4444.
- [12] S. Liu, Y. Cao, S. Kang, N. Hu, X. Liu, D. Su, D. Yu, and H. Meng, “Transferring Source Style in Non-Parallel Voice Conversion,” in *Proc. Interspeech*, 2020, pp. 4721–4725.
- [13] Q. Ma, R. Liu, X. Wen, C. Lu, and X. Chen, “Submission from SRCB for Voice Conversion Challenge 2020,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 131–135.
- [14] D. Stanton, Y. Wang, and RJ Skerry-Ryan, “Predicting Expressive Speaking Style from Text in End-To-End Speech Synthesis,” in *Proc. SLT*, 2018, pp. 595–602.
- [15] L. Sun, K. Li, H. Wang, S. Kang, and H. Meng, “Phonetic posteriorgrams for many-to-one voice conversion without parallel data training,” in *Proc. ICME*, 2016, pp. 1–6.
- [16] J. Zhang, Z. Ling, L. Liu, Y. Jiang, and L. Dai, “Sequence-to-Sequence Acoustic Modeling for Voice Conversion,” *IEEE/ACM TASLP*, vol. 27, no. 3, pp. 631–644, 2019.
- [17] W.-C. Huang, Y.-C. Wu, T. Hayashi, and T. Toda, “Any-to-One Sequence-to-Sequence Voice Conversion using Self-Supervised Discrete Speech Representations,” in *Proc. ICASSP*, 2021, pp. 5944–5948.
- [18] Y.-H. Peng, C.-H. Hu, A. Kang, H.-S. Lee, P.-Y. Chen, Y. Tsao, and H.-M. Wang, “The Academia Sinica Systems of Voice Conversion for VCC2020,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 180–183.
- [19] Y. Y. Lin, C.-M. Chien, J.-H. Lin, H.-Y. Lee, and L.-S. Lee, “FragmentVC: Any-to-Any Voice Conversion by End-to-End Extracting and Fusing Fine-Grained Voice Fragments With Attention,” in *Proc. ICASSP*, 2021, pp. 5939–5943.
- [20] A. Baevski, S. Schneider, and M. Auli, “vq-wav2vec: Self-Supervised Learning of Discrete Speech Representations,” in *Proc. ICLR*, 2020.
- [21] K. Inoue, S. Hara, M. Abe, T. Hayashi, R. Yamamoto, and S. Watanabe, “Semi-Supervised Speaker Adaptation for End-to-End Speech Synthesis with Pretrained Models,” in *Proc. ICASSP*, 2020, pp. 7634–7638.
- [22] L. Zheng, J. Tao, Z. Wen, and R. Zhong, “CASIA Voice Conversion System for the Voice Conversion Challenge 2020,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 136–139.
- [23] Z. Lian, R. Zhong, Z. Wen, B. Liu, and J. Tao, “Towards Fine-Grained Prosody Control for Voice Conversion,” in *Proc. ISCSLP*, 2021, pp. 1–5.
- [24] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “LibriSpeech: An ASR corpus based on public domain audio books,” in *Proc. ICASSP*, 2015, pp. 5206–5210.
- [25] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, “LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech,” in *Proc. Interspeech*, 2019, pp. 1526–1530.
- [26] K. Park and T. Mulc, “CSS10: A Collection of Single Speaker Speech Datasets for 10 Languages,” in *Proc. Interspeech*, 2019, pp. 1566–1570.
- [27] Munich Artificial Intelligence Laboratories GmbH, “The M-AIILABS speech dataset,” 2019, accessed 30 November 2019.
- [28] Data Baker China, “Chinese standard mandarin speech corpus,” accessed 05 May 2020.
- [29] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen, A. Renduchintala, and T. Ochiai, “ESPnet: End-to-End Speech Processing Toolkit,” in *Proc. Interspeech*, 2018, pp. 2207–2211.
- [30] T. Hayashi, R. Yamamoto, K. Inoue, T. Yoshimura, S. Watanabe, T. Toda, K. Takeda, Y. Zhang, and X. Tan, “Espnet-TTS: Unified, Reproducible, and Integratable Open Source End-to-End Text-to-Speech Toolkit,” in *Proc. ICASSP*, 2020, pp. 7654–7658.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N Gomez, L. Kaiser, and I. Polosukhin, “Attention is All you Need,” in *Proc. NIPS*, 2017, pp. 5998–6008.
- [32] L. Dong, S. Xu, and B. Xu, “Speech-Transformer: A No-Recurrence Sequence-to-Sequence Model for Speech Recognition,” in *Proc. ICASSP*, 2018, pp. 5884–5888.
- [33] S. Karita, N. E. Y. Soplin, S. Watanabe, M. Delcroix, A. Ogawa, and T. Nakatani, “Improving Transformer-Based End-to-End Speech Recognition with Connectionist Temporal Classification and Language Model Integration,” in *Proc. Interspeech*, 2019, pp. 1408–1412.
- [34] S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, “Hybrid CTC/Attention Architecture for End-to-End Speech Recognition,” *IEEE Jour. of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1240–1253, 2017.
- [35] D. Povey, G. Cheng, Y. Wang, K. Li, H. Xu, M. Yarmohammadi, and S. Khudanpur, “Semi-Orthogonal Low-Rank Matrix Factorization for Deep Neural Networks,” in *Proc. Interspeech*, 2018, pp. 3743–3747.

- [36] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, “fairseq: A Fast, Extensible Toolkit for Sequence Modeling,” in *Proc. NAACL*, 2019, pp. 48–53.
- [37] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, “X-vectors: Robust dnn embeddings for speaker recognition,” in *Proc. ICASSP*, 2018, pp. 5329–5333.
- [38] N. Li, S. Liu, Y. Liu, S. Zhao, and M. Liu, “Neural Speech Synthesis with Transformer Network,” in *Proc. AAAI*, 2019, pp. 6706–6713.
- [39] W.-C. Huang, T. Hayashi, Y.-C. Wu, H. Kameoka, and T. Toda, “Voice Transformer Network: Sequence-to-Sequence Voice Conversion Using Transformer with Text-to-Speech Pretraining,” in *Proc. Interspeech*, 2020, pp. 4676–4680.
- [40] W.-C. Huang, T. Hayashi, Y.-C. Wu, H. Kameoka, and T. Toda, “Pretraining Techniques for Sequence-to-Sequence Voice Conversion,” *IEEE TASLP*, vol. 29, pp. 745–755, 2021.
- [41] R. Yamamoto, E. Song, and J. Kim, “Parallel WaveGAN: A Fast Waveform Generation Model Based on Generative Adversarial Networks with Multi-Resolution Spectrogram,” in *Proc. ICASSP*, 2020, pp. 6199–6203.
- [42] R. Kumar Das, T. Kinnunen, W.-C. Huang, Z.-H. Ling, J. Yamagishi, Z. Yi, X. Tian, and T. Toda, “Predictions of Subjective Ratings and Spoofing Assessments of Voice Conversion Challenge 2020 Submissions,” in *Proc. Joint Workshop for the BC and VCC 2020*, 2020, pp. 99–120.
- [43] M. Morise, F. Yokomori, and K. Ozawa, “WORLD: A Vocoder-Based High-Quality Speech Synthesis System for Real-Time Applications,” *IEICE Transactions on Information and Systems*, vol. 99, pp. 1877–1884, 2016.
- [44] B. Naderi and R. Cutler, “An Open Source Implementation of ITU-T Recommendation P.808 with Validation,” in *Proc. Interspeech*, 2020, pp. 2862–2866.
- [45] ITU-T Recommendation P.808, “Subjective evaluation of speech quality with a crowdsourcing approach,” 2018.
- [46] L. van der Maaten and G. Hinton, “Visualizing data using t-SNE,” *JMLR*, vol. 9, pp. 2579–2605, 2008.
- [47] E. Battenberg, S. Mariooryad, D. Stanton, RJ Skerry-Ryan, M. Shannon, D. Kao, and T. Bagby, “Effective use of variational embedding capacity in expressive end-to-end speech synthesis,” *arXiv preprint arXiv:1906.03402*, 2019.