# Many-to-Many Voice Transformer Network

Hirokazu Kameoka, Wen-Chin Huang, Kou Tanaka,
Takuhiro Kaneko, Nobukatsu Hojo, and Tomoki Toda

*Abstract*—This paper proposes a voice conversion (VC) method based on a sequence-to-sequence (S2S) learning framework, which enables simultaneous conversion of the voice characteristics, pitch contour, and duration of input speech. We previously proposed an S2S-based VC method using a transformer network architecture called the voice transformer network (VTN). The original VTN was designed to learn only a mapping of speech feature sequences from one speaker to another. The main idea we propose is an extension of the original VTN that can simultaneously learn mappings among multiple speakers. This extension called the many-to-many VTN makes it able to fully use available training data collected from multiple speakers by capturing common latent features that can be shared across different speakers. It also allows us to introduce a training loss called the identity mapping loss to ensure that the input feature sequence will remain unchanged when the source and target speaker indices are the same. Using this particular loss for model training has been found to be extremely effective in improving the performance of the model at test time. We conducted speaker identity conversion experiments and found that our model obtained higher sound quality and speaker similarity than baseline methods. We also found that our model, with a slight modification to its architecture, could handle any-to-many conversion tasks reasonably well.

*Index Terms*—Voice conversion (VC), sequence-to-sequence learning, attention, transformer network, many-to-many VC.

## I. INTRODUCTION

Voice conversion (VC) is a technique to modify the voice characteristic and speaking style of an input utterance without changing the linguistic content. Potential applications of VC include speaker-identity modification [1], speaking aids [2], [3], speech enhancement [4], [5], and accent conversion [6].

Many conventional VC methods use parallel utterances of source and target speech to train acoustic models for feature mapping. The training process typically consists of computing acoustic features from source and target utterances, using dynamic time warping (DTW) to align parallel utterances, and training the acoustic model that describes the mappings between the source features and the target features. Examples of acoustic models include Gaussian mixture models (GMMs) [7], [8], partial least square regression [9], dynamic kernel partial least square regression [10], frequency warping [11], non-negative matrix factorization [12], group sparse representation [13], fully-connected deep neural networks (DNNs) [14], [15], and long-short term memory networks (LSTMs) [16], [17], just to name a few. Recently, attempts have also been made to

H. Kameoka, K. Tanaka, T. Kaneko and N. Hojo are with NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation, Atsugi, Kanagawa, 243-0198 Japan (e-mail: hirokazu.kameoka.uh@hco.ntt.co.jp) and W.-C. Huang and T. Toda are with Nagoya University.

formulate non-parallel methods using deep generative models. These include methods based on variational autoencoders [18]–[22], generative adversarial networks [23], [24], and flow-based models [25]. Interested readers are referred to an excellent review article [26] on the recent trends in VC research.

Although most of the methods mentioned above are successful in converting the local spectral features, they have difficulty in converting suprasegmental features that reflect long-term dependencies, such as the fundamental frequency ($F_0$) contour, duration and rhythm of the input speech. This is because acoustic models in these methods can only describe mappings between local features. However, since these suprasegmental features are as important factors as local features that characterize speaker identities and speaking styles, it would be desirable if these features could also be converted more flexibly. Although some attempts have been made to achieve conversions of $F_0$ and energy contours through their continuous wavelet transform representations [17], [27]–[30], they do not address conversions of duration and rhythm. One solution to overcome this limitation would be to adopt a sequence-to-sequence (seq2seq or S2S) model, as it has a powerful ability to learn mappings between sequential data of variable lengths by capturing and using long-range dependencies.

S2S models [31], [32] have recently been applied with notable success in many tasks such as machine translation, automatic speech recognition (ASR) [32] and text-to-speech (TTS) [33]–[38]. They are composed mainly of two elements: an encoder and decoder. The encoder encodes an input sequence to its latent representation in the form of hidden state vectors whereas the decoder generates an output sequence according to this latent representation. With the original S2S model, all input sequences are encoded into a single context vector of a fixed dimension. One problem with this is that the ability of the model to capture long-range dependencies can be limited especially when input sequences are long. To overcome this limitation, a mechanism called attention [39] has been introduced, which allows the network to learn where to pay attention in the input sequence when producing each item in the output sequence.

While recurrent neural networks (RNNs) have initially been used as the default option for designing the encoder and decoder networks in S2S models, recent work has shown that convolutional neural network (CNN)-based architectures also have excellent potential for capturing long-term dependencies [40]. Subsequently, yet another type of architecture called the transformer has been proposed [41], which uses neither convolution nor recurrent layers in its network, but only the mechanism of attention. In particular, it uses multi-head self-attention layers to design the two networks. Self-attention is a

type of attention mechanism, which offers an efficient way to relate different positions of a given sequence. The multi-head self-attention mechanism splits each vector in a sequence into smaller parts and then computes the self-attention over the sequence of each part in parallel. Unlike RNNs, both these architectures have the advantage that they are suitable for parallel computations using GPUs.

Several VC methods based on S2S models have already been proposed [42]–[45]. We also previously proposed VC methods based on S2S models with RNN [46], CNN [47] and transformer architectures [48]. In our most recent work [48], we proposed a VC method based on a transformer architecture, which we called the voice transformer network (VTN). Through this work, it transpired that the model trained from scratch did not perform as expected when the amount of training data was limited. To address this, we introduced a TTS pretraining technique to provide a good initialization for fast and sample-efficient VC model training with the aid of text-speech pair data, thus reducing the parallel data size requirement and training time.

In this paper, we propose several ideas to make the VTN perform well even without TTS pretraining. One limitation with regular S2S models including the VTN is that they can only learn a mapping from one speaker to another. When parallel utterances of multiple speakers are available, naively preparing and training a different model independently for each speaker pair would be inefficient, since the model for a particular speaker pair fails to use the training data of the other speakers for its training. To fully utilize available training data in multiple speakers, the main idea we propose is extending the VTN so that it can simultaneously learn mappings among multiple speakers. We call this extended version the many-to-many VTN. The idea of this many-to-many extension was triggered by our previous studies on non-parallel VC methods [49] that showed the significant superiority of a many-to-many VC method [24] over its one-to-one counterpart [23]. This suggests the possibility that in learning mappings between a certain speaker pair the data of other speakers may also be useful. However, the effect this extension would have on parallel VC tasks based on S2S models was nontrivial. In this regard, we consider the present contribution important. We further propose several ideas, including the identity mapping loss, attention windowing, and settings needed to achieve any-to-many and real-time conversions. We show that the performance improvement achieved by the many-to-many extension owes not only to the fact that it can leverage more training data than the one-to-one counterpart but also to the identity mapping loss, which is only available in the many-to-many version. We also show that the pre-layer normalization (Pre-LN) architecture, proposed in [50], [51], is effective in both the original (one-to-one) and many-to-many VTNs.

## II. RELATED WORK

Several VC methods based on S2S models have already been proposed, including the ones we proposed previously. Regular S2S models usually require large-scale parallel corpora for training. However, collecting parallel utterances is often costly and non-scalable. Therefore, in VC tasks using S2S models, one challenge is how to make them work with a limited amount of training data.

One solution is to use text labels as auxiliary information for model training. Miyoshi et al. proposed a VC model combining an S2S model and acoustic models for ASR and TTS [42]. Zhang et al. proposed an S2S-based VC model guided by an ASR system [43]. Subsequently, Zhang et al. proposed a shared S2S model for TTS and VC tasks [44]. Recently, Biadsy et al. proposed an end-to-end VC model called "Parrotron", which uses a multitask learning strategy to train the encoder and decoder along with an ASR model [45]. Our VTN [48] is another example, which relies on TTS pretraining using text-speech pair data.

The proposed many-to-many VTN differs from the above methods in that it does not rely on ASR or TTS models and requires no text annotations for model training.

## III. VOICE TRANSFORMER NETWORK

### A. Feature extraction and normalization

Following our previous work [47], in this work, we choose to use the mel-cepstral coefficients (MCCs) [52], log $F_0$, aperiodicity, and voiced/unvoiced indicator of speech as the acoustic features to be converted. Once acoustic features have successfully been converted, we can use either a conventional vocoder or one of the recently developed high-quality neural vocoders [53]–[65] to generate the signals of converted speech.

To obtain an acoustic feature vector, we follow the same procedure in [47]. Namely, we first use WORLD [66] to analyze the spectral envelope, log $F_0$, coded aperiodicity, and voiced/unvoiced indicator within each time frame of a speech utterance, then compute $I$ MCCs from the extracted spectral envelope, and finally construct an acoustic feature vector by stacking the MCCs, log $F_0$, coded aperiodicity, and voiced/unvoiced indicator. Each acoustic feature vector thus consists of $I + 3$ elements. At training time, we normalize each element $x_{i,n}$ $(i = 1, \ldots, I)$ of the MCCs and log $F_0$ $x_{I+1,n}$ at frame $n$ to $x_{i,n} \leftarrow (x_{i,n} - \mu_i)/\sigma_i$ where $i$, $\mu_i$, and $\sigma_i$ denote the feature index, mean, and standard deviation of the $i$-th feature within all the voiced segments of the training samples of the same speaker.

As with our previous work [47], we found it useful to use a similar trick introduced in previous work [67] to accelerate and stabilize training and inference. Namely, we divide the acoustic feature sequence obtained above into non-overlapping segments of equal length $r$ and use the stack of the acoustic feature vectors in each segment as a new feature vector so that the new feature sequence becomes $r$ times shorter than the original feature sequence.

### B. Model

We hereafter use $\mathbf{X}^{(\mathsf{s})} = [\mathbf{x}_1^{(\mathsf{s})}, \ldots, \mathbf{x}_{N_\mathsf{s}}^{(\mathsf{s})}] \in \mathbb{R}^{D \times N_\mathsf{s}}$ and $\mathbf{X}^{(\mathsf{t})} = [\mathbf{x}_1^{(\mathsf{t})}, \ldots, \mathbf{x}_{N_\mathsf{t}}^{(\mathsf{t})}] \in \mathbb{R}^{D \times N_\mathsf{t}}$ to denote the source and target speech feature sequences of non-aligned parallel utterances, where $N_\mathsf{s}$ and $N_\mathsf{t}$ denote the lengths of the two sequences and $D$ denotes the feature dimension. The VTN
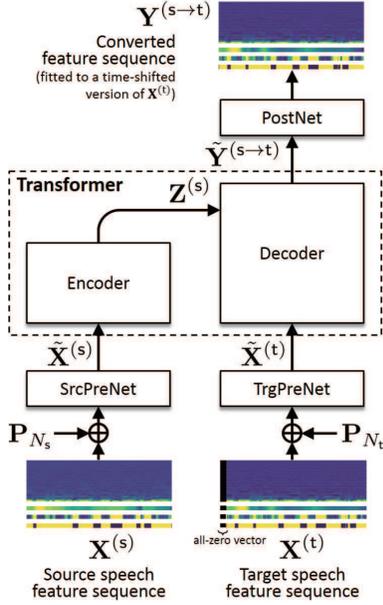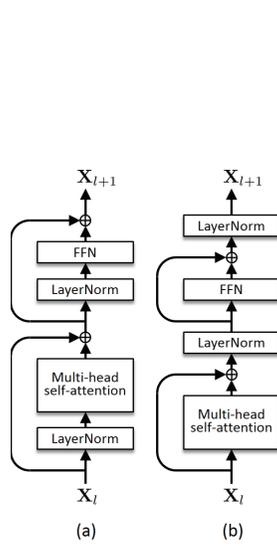
Fig. 1. Overall structure of VTN



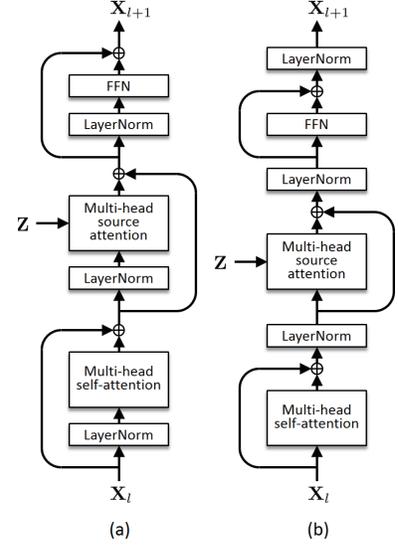Fig. 2. Encoder layers with (a) Pre-LN and (b) Post-LN architectures



Fig. 3. Decoder layers with (a) Pre-LN and (b) Post-LN architectures

[48] has an encoder-decoder structure with a transformer architecture [41] that maps $\mathbf{X}^{(s)}$ to $\mathbf{X}^{(t)}$ (Fig. 1). The encoder is expected to extract contextual information from source speech, and the decoder produces the target speech feature sequence according to the contextual information the encoder has generated. Unlike RNN- and CNN-based S2S models, the transformer model does not have any sense of the order of the elements in a sequence. Thus, the sinusoidal position encodings [41], denoted by $\mathbf{P}_{N_s} \in \mathbb{R}^{D \times N_s}$ and $\mathbf{P}_{N_t} \in \mathbb{R}^{D \times N_t}$, are first added to the source and target feature sequences to make the model "aware" of the position at which an each element in the sequences is located. The source and target feature sequences are then passed through convolutional prenets, which we call the source and target prenets, before being fed into the encoder and decoder. The output $\tilde{\mathbf{Y}}^{(s \to t)}$ from the decoder is finally passed through a convolutional postnet before producing the final output $\mathbf{Y}^{(s \to t)}$. The two prenets and the postnet, each consisting of three convolution layers, are used to capture and express the local dynamics in source and target speech and convert input sequences into sequences of the same lengths. In the following, we use $\tilde{\mathbf{X}}^{(s)} \in \mathbb{R}^{d \times N_s}$ and $\tilde{\mathbf{X}}^{(t)} \in \mathbb{R}^{d \times N_t}$ to denote the outputs from the source and target prenets, respectively, where $d$ is the output channel number of each prenet.

*1) Encoder:* The encoder takes $\tilde{\mathbf{X}}^{(s)}$ as the input and produces a context vector sequence $\mathbf{Z}^{(s)} = [\mathbf{z}_1^{(s)}, \ldots, \mathbf{z}_{N_s}^{(s)}] \in \mathbb{R}^{d \times N_s}$. The encoder consists of $L$ identical layers, each of which has self-attention (SA) and position-wise fully-connected feed forward network (FFN) sub-layers. Residual connections and layer normalizations are applied in addition to the two sub-layers.

**Multi-head self-attention sub-layer:** By using $\mathbf{X} \in \mathbb{R}^{d \times N}$ and $\mathbf{Y} \in \mathbb{R}^{d \times N}$ to denote the input and output sequences of length $N$ of an SA sub-layer, the process $\mathbf{Y} = \mathrm{SA}(\mathbf{X})$, by

which $\mathbf{Y}$ is produced, is given as

$$[\mathbf{Q}; \mathbf{K}; \mathbf{V}] = \mathbf{W}_1 \mathbf{X} \in \mathbb{R}^{3d \times N}, \tag{1}$$

$$\text{where} \begin{cases} \mathbf{Q} = [\mathbf{Q}_1; \ldots; \mathbf{Q}_H] \\ \mathbf{K} = [\mathbf{K}_1; \ldots; \mathbf{K}_H] \\ \mathbf{V} = [\mathbf{V}_1; \ldots; \mathbf{V}_H] \end{cases}, \tag{2}$$

$$\mathbf{A}_h = \mathsf{softmax}\left(\frac{\mathbf{K}_h^\mathsf{T} \mathbf{Q}_h}{\sqrt{d}}\right) \ (h = 1, \ldots, H), \tag{3}$$

$$\mathbf{Y} = \mathbf{W}_2 [\mathbf{V}_1 \mathbf{A}_1; \ldots; \mathbf{V}_H \mathbf{A}_H], \tag{4}$$

where $\mathbf{W}_1 \in \mathbb{R}^{3d \times d}$ and $\mathbf{W}_2 \in \mathbb{R}^{d \times d}$ are learnable weight matrices, $\mathsf{softmax}$ denotes a softmax operation performed on the first axis, $H$ denotes the number of heads, and $[; ]$ denotes vertical concatenation of matrices (or vectors) with compatible sizes. Intuitively, this process can be understood as follows. First, an input vector sequence is converted into three types of vector sequences with the same shape, which can be metaphorically interpreted as the queries and the key-value pairs in a hash table. Each of the three vector sequences is further split into $H$ homogeneous vector sequences with the same shape. By using the query and key pair, Eq. (3) computes a self-attention matrix, whose element measures how contextually similar each pair of vectors is in the given sequence $\mathbf{X}$. The splitting into $H$ heads allows us to measure self-simiarity in terms of $H$ different types of context. The $n$-th column of $\mathbf{V}_h \mathbf{A}_h$ in Eq. (4) can be seen as a new feature vector given by activating the value vectors at all the positions that are similar to the current position $n$ in terms of context $h$ and adding them together. Eq. (4) finally produces the output sequence $\mathbf{Y}$ after combining all these feature vector sequences using learnable weights.

**Position-wise feed forward network sub-layer:** By using $\mathbf{X} \in \mathbb{R}^{d \times N}$ and $\mathbf{Y} \in \mathbb{R}^{d \times N}$ again to denote the input and output sequences of length $N$ of an FFN sub-layer, the process

$\mathbf{Y} = \mathsf{FFN}(\mathbf{X})$, by which $\mathbf{Y}$ is produced, is given as

$$\mathbf{Y} = \mathbf{W}_4 \phi(\mathbf{W}_3 \mathbf{X} + \mathbf{B}_3) + \mathbf{B}_4, \tag{5}$$

where $\mathbf{W}_3 \in \mathbb{R}^{d' \times d}$, $\mathbf{W}_4 \in \mathbb{R}^{d \times d'}$ are learnable weight matrices, $\mathbf{B}_3 = [\mathbf{b}_3, \dots, \mathbf{b}_3] \in \mathbb{R}^{d' \times N}$ and $\mathbf{B}_4 = [\mathbf{b}_4, \dots, \mathbf{b}_4] \in \mathbb{R}^{d \times N}$ are bias matrices, each consisting of identical learnable column vectors, and $\phi$ denotes an elementwise nonlinear activation function such as the rectified linear unit (ReLU) and gated linear unit (GLU) functions.

**Layer normalization sub-layers:** Recent work has shown that the location of the layer normalization in the transformer architecture affects the speed and stability of the training process as well as the performance of the trained model [50], [51]. While the original transformer architecture places layer normalization after the SA and FFN sub-layers, the architectures presented in the subsequent work [50], [51] are designed to place it before them, as illustrated in Fig. 2. To distinguish between these two architectures, we refer to the former and latter as post-layer normalization (Post-LN) and Pre-LN architectures, respectively. We show later how differently these architectures actually performed in our experiments. Note that when we say we apply layer normalization to an input vector sequence, say $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$, we mean applying layer normalization to all the vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ treated as mini-batch samples.

If we use $\mathbf{X}_l$ and $\mathbf{X}_{l+1}$ to denote the input and output of the $l$-th encoder layer (with the PreLN architecture), the process $\mathbf{X}_{l+1} = \mathsf{Enc}_l(\mathbf{X}_l)$ of the $l$-the layer is given by

$$\mathbf{U} = \mathbf{X}_l + \mathsf{SA}(\mathsf{LayerNorm}_1(\mathbf{X}_l)), \tag{6}$$

$$\mathbf{X}_{l+1} = \mathbf{U} + \mathsf{FFN}(\mathsf{LayerNorm}_2(\mathbf{U})), \tag{7}$$

where $\mathsf{LayerNorm}_1$ and $\mathsf{LayerNorm}_2$ denote different LN sub-layers. As described above, each layer has learnable parameters in the SA and FFN sub-layers and the two LN sub-layers. The layer implemented as above is particularly attractive in that it is able to relate all the positions in the entire input sequence using only a single layer. This is in contrast to a regular convolution layer, which is only able to relate local positions near each position.

*2) Decoder:* The decoder takes $\mathbf{Z}^{(\mathsf{s})}$ and $\tilde{\mathbf{X}}^{(\mathsf{t})}$ as the inputs and produces a converted feature sequence $\tilde{\mathbf{Y}}^{(\mathsf{s} \to \mathsf{t})} = [\tilde{\mathbf{y}}_1^{(\mathsf{s} \to \mathsf{t})}, \dots, \tilde{\mathbf{y}}_{N_t}^{(\mathsf{s} \to \mathsf{t})}] \in \mathbb{R}^{d \times N_t}$. Similar to the encoder, the decoder consists of $L$ identical layers, each of which has SA and FFN sub-layers, residual connections, and layer normalization sub-layers. In addition to these sub-layers, each layer has a multi-head target-to-source attention (TSA) sub-layer, as illustrated in Fig. 3, whose role is to find which position in the source feature sequence contextually corresponds to each position in the target feature sequence and convert the context vector sequence according to the predicted corresponding positions.

**Multi-head target-to-source attention sub-layer:** By using $\mathbf{X} \in \mathbb{R}^{d \times N}$ and $\mathbf{Y} \in \mathbb{R}^{d \times N}$ to denote the input and output sequences of length $N$ of the TSA sub-layer, the process $\mathbf{Y} = \mathsf{TSA}(\mathbf{X}, \mathbf{Z})$, by which $\mathbf{Y}$ is produced, is given in the same manner as the SA sub-layer with the only difference being that the key and value pair $(\mathbf{K}, \mathbf{V})$ is computed using the

output $\mathbf{Z}$ from the encoder:

$$\mathbf{Q} = \mathbf{W}_5 \mathbf{X}, \tag{8}$$

$$[\mathbf{K}; \mathbf{V}] = \mathbf{W}_6 \mathbf{Z}, \tag{9}$$

$$\text{where} \begin{cases} \mathbf{Q} = [\mathbf{Q}_1; \dots; \mathbf{Q}_H] \\ \mathbf{K} = [\mathbf{K}_1; \dots; \mathbf{K}_H] \\ \mathbf{V} = [\mathbf{V}_1; \dots; \mathbf{V}_H] \end{cases}, \tag{10}$$

$$\mathbf{A}_h = \mathsf{softmax}\left(\frac{\mathbf{K}_h^\mathsf{T} \mathbf{Q}_h}{\sqrt{d}}\right) \ (h = 1, \dots, H), \tag{11}$$

$$\mathbf{Y} = \mathbf{W}_7 [\mathbf{V}_1 \mathbf{A}_1; \dots; \mathbf{V}_H \mathbf{A}_H], \tag{12}$$

where $\mathbf{W}_5 \in \mathbb{R}^{d \times d}$, $\mathbf{W}_6 \in \mathbb{R}^{2d \times d}$, and $\mathbf{W}_7 \in \mathbb{R}^{d \times d}$ are learnable weight matrices. Analogously to the SA sub-layer, Eq. (11) is used to compute a TSA matrix using the query and key pair, where the $(n, m)$-th element indicates the similarity between the $n$-th and $m$-th frames of source and target speech. The peak trajectory of $\mathbf{A}_h$ can thus be interpreted as a time warping function that associates the frames of the source speech with those of the target speech. The splitting into $H$ heads allows us to measure the simiarity in terms of $H$ different types of context. $\mathbf{V}_h \mathbf{A}_h$ in Eq. (12) can be thought of as a time-warped version of $\mathbf{V}_h$ in terms of context $h$. Eq. (12) finally produces the output sequence $\mathbf{Y}$ after combining all these time-warped feature sequences using learnable weights.

All the other sub-layers are defined in the same way as the encoder. The overall structures of the decoder layers with the PreLN and PostLN architectures are depicted in Fig. 3. If we use $\mathbf{X}_l$ and $\mathbf{X}_{l+1}$ to denote the input and output of the $l$-th decoder layer (with the PreLN architecture), the process $\mathbf{X}_{l+1} = \mathsf{Dec}(\mathbf{X}_l, \mathbf{Z})$ of the $l$-th layer is given by

$$\mathbf{U}_1 = \mathbf{X}_l + \mathsf{SA}(\mathsf{LayerNorm}_1(\mathbf{X}_l)), \tag{13}$$

$$\mathbf{U}_2 = \mathbf{U}_1 + \mathsf{TSA}(\mathsf{LayerNorm}_2(\mathbf{U}_1), \mathbf{Z}), \tag{14}$$

$$\mathbf{X}_{l+1} = \mathbf{U}_2 + \mathsf{FFN}(\mathsf{LayerNorm}_3(\mathbf{U}_2)). \tag{15}$$

Note that each layer has learnable parameters in the SA, FFN, and TSA sub-layers and the three LN sub-layers.

*3) Autoregressive structure:* Since the target feature sequence $\mathbf{X}^{(\mathsf{t})}$ is of course not accessible at test time, we would want to use a feature vector that the decoder has generated as the input to the decoder for the next time step so that feature vectors can be generated one-by-one in a recursive manner. To allow the model to behave in this manner, we must first take care that the decoder must not be allowed to use future information about the target feature vectors when producing an output vector at each time step. This can be ensured by simply constraining the convolution layers in the target prenet to be causal and replacing Eq. (3) in all the SA sub-layers in the decoder with

$$\mathbf{A}_h = \mathsf{softmax}\left(\frac{\mathbf{K}_h^\mathsf{T} \mathbf{Q}_h}{\sqrt{d}} + \mathbf{E}\right), \tag{16}$$

where $\mathbf{E}$ is a matrix whose $(n, n')$-th element is given by

$$e_{n,n'} = \begin{cases} 0 & (n \le n') \\ -\infty & (n > n') \end{cases}, \tag{17}$$

so that the predictions for position $n$ can depend only on the known outputs at positions less than $n$. Second, an all-zero vector is appended to the left end of $\mathbf{X}^{(\mathsf{t})}$

$$\mathbf{X}^{(\mathsf{t})} \leftarrow [\mathbf{0}, \mathbf{X}^{(\mathsf{t})}], \tag{18}$$

so that the recursion always begins with the all-zero vector. Third, the output sequence $\mathbf{Y}^{(s \to t)}$ must correspond to a time-shifted version of $\mathbf{X}^{(t)}$ so that at each time step the decoder will be able to predict the target speech feature vector that is likely to appear at the next time step. To this end, we include an $L_1$ loss

$$\mathcal{L}_{\text{main}} = \tfrac{1}{M}\|[\mathbf{Y}^{(s \to t)}]_{:,1:M} - [\mathbf{X}^{(t)}]_{:,2:M+1}\|_1, \quad (19)$$

in the training loss to be minimized, where we have used the colon operator : to specify the range of indices of the elements in a matrix we wish to extract. For ease of notation, we use : to represent all elements along an axis.

### C. Constraints on Attention Matrix

Since the alignment path between parallel utterances must lie close to the diagonal, the diagonal region in the attention matrices in each TSA sub-layer in the decoder should be dominant. By imposing such a restriction, the search space during training can be significantly reduced, thus significantly reducing the training effort. One way to force the attention matrices to be diagonally dominant involves introducing a diagonal attention loss (DAL) [37]:

$$\mathcal{L}_{\text{dal}} = \tfrac{1}{NMLH} \sum_l \sum_h \|\mathbf{G}_{N_s \times N_t} \odot \mathbf{A}_{l,h}\|_1, \quad (20)$$

where $\mathbf{A}_{l,h}$ denotes the target-to-source attention matrix of the $h$-th head in the TSA sub-layer in the $l$-th decoder layer, $\odot$ denotes an elementwise product, and $\mathbf{G}_{N_s \times N_t} \in \mathbb{R}^{N_s \times N_t}$ is a non-negative weight matrix whose $(n, m)$-th element $w_{n,m}$ is defined as $w_{n,m} = 1 - e^{-(n/N_s - m/N_t)^2/2\nu^2}$.

### D. Training loss

Given examples of parallel utterances, the total training loss for the VTN to be minimized is given as

$$\mathcal{L} = \mathbb{E}_{\mathbf{X}^{(s)}, \mathbf{X}^{(t)}} \{\mathcal{L}_{\text{main}} + \lambda_{\text{dal}} \mathcal{L}_{\text{dal}}\}, \quad (21)$$

where $\mathbb{E}_{\mathbf{X}^{(s)}, \mathbf{X}^{(t)}}\{\cdot\}$ denotes the sample mean over all the training examples and $\lambda_{\text{dal}} \geq 0$ is a regularization parameter, which weighs the importance of $\mathcal{L}_{\text{dal}}$ relative to $\mathcal{L}_{\text{dec}}$.

### E. Conversion Algorithm

At test time, a source speech feature sequence $\mathbf{X}$ can be converted via Algorithm 1. Note that in our model, the all-zero vector corresponds to the start-of-sequence token. As for the end-of-sequence token, we intentionally did not include it in the source and target feature sequences. This is because we are assuming a situation where source speech features are constantly coming in and the conversion is performed online. In the following experiments, we set $M$ to a sufficiently large number (specifically, twice the length $N$ of the source feature sequence) and regarded the time $m$ at which the attended time point (i.e., the peak of the attention distribution) first reached $N$ as the end-of-utterance.

Once $\mathbf{Y}$ has been obtained, we adjust the mean and variance of the generated feature sequence so that they match the pretrained mean and variance of the feature vectors of the target speaker. We can then generate a time-domain signal

---

**Algorithm 1** Default conversion algorithm

$\mathbf{Z} \leftarrow \mathbf{X}, \mathbf{Y} \leftarrow \mathbf{0}$
**for** $l = 1$ to $L$ **do**
  $\mathbf{Z} \leftarrow \text{Enc}_l(\mathbf{Z})$
**end for**
**for** $m = 1$ to $M$ **do**
  **for** $l = 1$ to $L$ **do**
    $\mathbf{Y} \leftarrow \text{Dec}_l(\mathbf{Y}, \mathbf{Z})$
  **end for**
  $\mathbf{Y} \leftarrow [\mathbf{0}, \mathbf{Y}]$
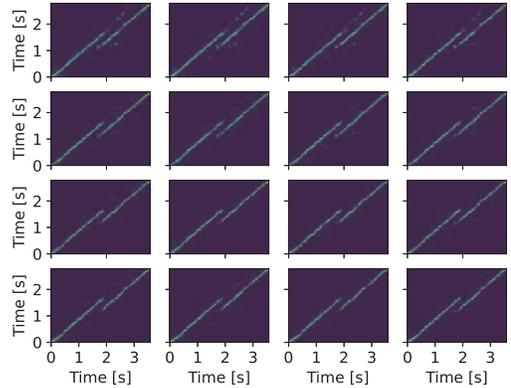**end for**
**return** $\mathbf{Y}$

---



Fig. 4. Example of TSA matrices predicted using the original VTN with $L = 4$ and $H = 4$ trained from scratch (without pretraining). Graph of column $h$ and row $l$ shows plot of $\mathbf{A}_{l,h}$.

using the WORLD vocoder or any recently developed neural vocoder.

However, as Fig. 4 shows, it transpired that with the model trained from scratch, the attended time point did not always move forward monotonically and continuously at test time and occasionally made a sudden jump to a distant time point, resulting in some segments being skipped or repeated, even though the DAL was considered in training. In our previous work [48], we proposed to introduce pretraining techniques exploiting auxiliary text labels to improve the behavior and performance of the conversion algorithm, as mentioned earlier. In the next section, we present several ideas that can greatly improve the behavior of the VTN even without pretraining using text labels.

## IV. MANY-TO-MANY VTN

### A. Many-to-Many Extension

The first and main idea is a many-to-many extension of the VTN, which uses a single model to enable mappings among multiple speakers by allowing the prenets, postnet, encoder, and decoder to take source and target speaker indices as additional inputs. The overall structure of the many-to-many VTN is shown in Fig. 5.

Let $\mathbf{X}^{(1)}, \ldots, \mathbf{X}^{(K)}$ be examples of the acoustic feature sequences of different speakers reading the same sentence. Given a single pair of parallel utterances $\mathbf{X}^{(k)}$ and $\mathbf{X}^{(k')}$, where $k$ and $k'$ denote the source and target speaker indices
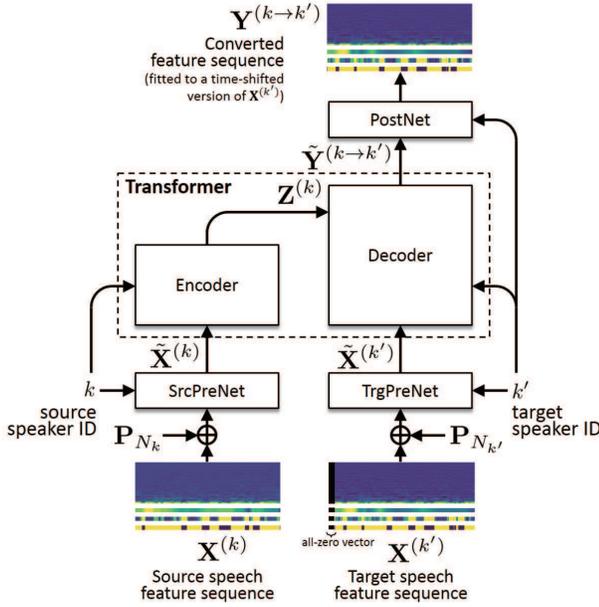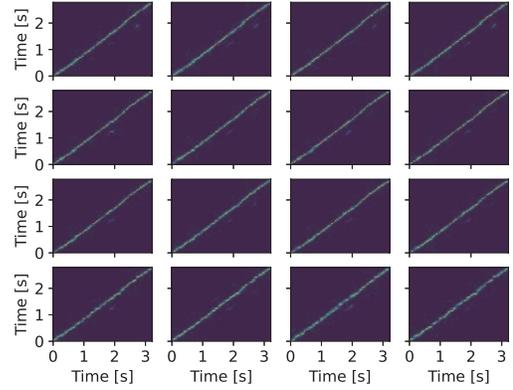
Fig. 5. Structure of many-to-many VTN.



Fig. 6. Example of attention matrices predicted using many-to-many VTN with $L = 4$ and $H = 4$ trained from scratch. Graph of column $h$ and row $l$ shows plot of $\mathbf{A}_{l,h}$.
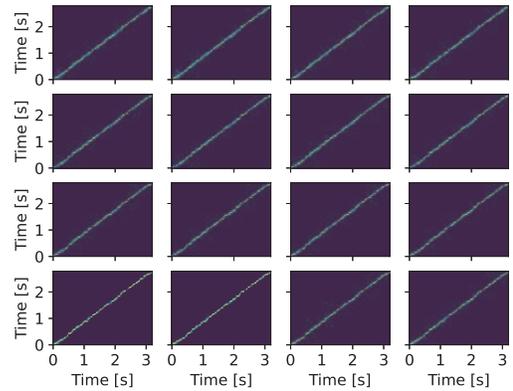


Fig. 7. Example of attention matrices predicted using the windowing technique based on original VTN with $L = 4$ and $H = 4$ trained from scratch. Graph of column $h$ and row $l$ shows plot of $\mathbf{A}_{l,h}$.

(integers), the source and target prenets take tuples $(\mathbf{X}^{(k)}, k)$ and $(\mathbf{X}^{(k')}, k')$ as the inputs and produce modified feature sequences $\tilde{\mathbf{X}}^{(k)}$ and $\tilde{\mathbf{X}}^{(k')}$, respectively. The encoder takes a tuple $(\tilde{\mathbf{X}}^{(k)}, k)$ as the input and produces a context vector sequence $\mathbf{Z}^{(k)}$. The decoder takes $(\tilde{\mathbf{X}}^{(k')}, \mathbf{Z}^{(k)}, k')$ as the input and produces a converted feature sequence $\tilde{\mathbf{Y}}^{(k \to k')}$. The postnet takes $(\tilde{\mathbf{Y}}^{(k \to k')}, k')$ as the input and finally produces a modified version $\mathbf{Y}^{(k \to k')}$ of $\tilde{\mathbf{Y}}^{(k \to k')}$. Each of the networks incorporates the speaker index into its process by modifying the input sequence, say $\mathbf{X}$, via

$$\mathbf{S} = \text{repeat}(\text{embed}(k)), \tag{22}$$
$$\mathbf{X} \leftarrow [\mathbf{X}; \mathbf{S}], \tag{23}$$

every time before feeding $\mathbf{X}$ into the SA, FFN, or TSA sublayers, where embed denotes an operation that retrieves a continuous vector from an integer input and repeat denotes an operation that produces a vector sequence from an input vector by simply repeating it along the time axis. Again, we append an all-zero vector to the left end of $\mathbf{X}^{(k')}$

$$\mathbf{X}^{(k')} \leftarrow [\mathbf{0}, \mathbf{X}^{(k')}]. \tag{24}$$

The loss functions to be minimized given this single training example are given as

$$\mathcal{L}_{\text{main}}^{(k,k')} = \frac{1}{N_{k'}} \| [\mathbf{Y}^{(k \to k')}]_{:,1:N_{k'}} - [\mathbf{X}^{(k')}]_{:,2:N_{k'}+1} \|_1, \tag{25}$$
$$\mathcal{L}_{\text{dal}}^{(k,k')} = \frac{1}{N_k N_{k'} HL} \sum_h \sum_l \| \mathbf{G}_{N_k \times N_{k'}} \odot \mathbf{A}_{l,h}^{(k,k')} \|_1, \tag{26}$$

where $\mathbf{A}_{l,h}^{(k,k')}$ denotes the TSA matrix of the $h$-th head in the TSA sub-layer in the $l$-th decoder layer.

With the above model, we can also consider the case in which $k = k'$. Minimizing the sum of the above losses under $k = k'$ encourages the model to let the input feature sequence $\mathbf{X}^{(k)}$ remain unchanged when $k$ and $k'$ indicate the same speaker. We call this loss the identity mapping loss (IML). The task of reconstructing an input sequence is

relatively easier than that of finding the mapping from an input sequence to a different sequence. Using the IML as the training objective allows the model to concentrate on learning only the autoregressive behaviour of the decoder. This can contribute to making the model training easier. The total training loss including the IML thus becomes

$$\mathcal{L} = \sum_{k,k' \neq k} \mathbb{E}_{\mathbf{X}^{(k)}, \mathbf{X}^{(k')}} \{ \mathcal{L}_{\text{all}}^{(k,k')} \} + \lambda_{\text{iml}} \sum_k \mathbb{E}_{\mathbf{X}^{(k)}} \{ \mathcal{L}_{\text{all}}^{(k,k)} \},$$
$$\text{where} \quad \mathcal{L}_{\text{all}}^{(k,k')} = \mathcal{L}_{\text{main}}^{(k,k')} + \lambda_{\text{dal}} \mathcal{L}_{\text{dal}}^{(k,k')}, \tag{27}$$

where $\mathbb{E}_{\mathbf{X}^{(k)}, \mathbf{X}^{(k')}} \{ \cdot \}$ and $\mathbb{E}_{\mathbf{X}^{(k)}} \{ \cdot \}$ denote the sample means over all the training examples of parallel utterances of speakers $k$ and $k'$, and $\lambda_{\text{iml}} \geq 0$ is a regularization parameter, which weighs the importance of the IML. The significant effect of the IML is shown later.

Fig. 6 shows examples of the TSA matrices predicted using the many-to-many VTN from the same test samples used in Fig. 4. As these examples show, the predicted attention matrices obtained with the many-to-many VTN exhibit more monotonic and continuous trajectories than those with the original VTN, demonstrating the impact of the many-to-many extension.

## B. Attention Windowing

We present another idea that can be used alone or combined with the many-to-many extension to improve the original VTN. To assist the attended point to move forward monotonically and continuously at test time, we propose to modify Algorithm 1 by adopting an idea inspired by the technique called *windowing* [32]. Specifically, we limit the paths through which the attended point is allowed to move by forcing the attentions to all the time points distant from the previous attended time point to zeros. We assume the attended time point to be the peak of the attention distribution, given as the mean of all the TSA attention matrices in the TSA sub-layers in the decoder. This can be implemented by replacing Eq. (11) in the TSA sub-layer in each decoder layer $l$ at the $m'$-th iteration of the for-loop for $m = 1, \ldots, M$ in the conversion process with

$$\hat{\mathbf{A}}_{l,h} = \mathsf{softmax}\left(\frac{\mathbf{K}_h^\mathsf{T}\mathbf{Q}_h}{\sqrt{d}} + \mathbf{F}\right) \ (h = 1, \ldots, H), \quad (28)$$

where the $(n,m)$-th element $f_{n,m}$ of $\mathbf{F}$ is given by

$$f_{n,m} = \begin{cases} -\infty & (m = m', \ n = 1, \ldots, \hat{n} - N_0) \\ -\infty & (m = m', \ n = \hat{n} + N_1, \ldots, N) , \\ 0 & (\text{otherwise}) \end{cases} \quad (29)$$

so that all the elements of the last column of the resulting $\hat{\mathbf{A}}_{l,h}$ become zero except for the elements from row $\max(1, \hat{n} - N_0)$ to row $\min(\hat{n} + N_1, N)$. $\mathbf{Z}$ denotes the final output of the encoder, $\mathbf{X}$ and $\mathbf{Y}$ denote the outputs of the previous and current sub-layers in the $l$-th decoder layer, and $\hat{n}$ denotes the maximum point of the attention distribution obtained at the $(m'-1)$-th iteration:

$$\hat{n} = \begin{cases} 1 & (m' = 1) \\ \mathrm{argmax}_n \frac{1}{LH} \sum_l \sum_h [\hat{\mathbf{A}}_{l,h}]_{:,m'-1} & (m' \neq 1) \end{cases} . \quad (30)$$

Note that we set $N_0$ and $N_1$ at the nearest integers that correspond to 160 and 320 ms, respectively. Fig. 7 shows examples of the TSA matrices obtained with this algorithm from the same test samples used in Fig. 4. As these examples show, this algorithm was found to have a certain effect on generating reasonably monotonic and continuous attention trajectories even without any modifications to the model structure of the original VTN.

It should be noted that we can also use the above algorithm as well as the algorithm presented in Subsection III-E for the many-to-many VTN, simply by replacing $\mathsf{Enc}_l(\mathbf{Z})$ and $\mathsf{Dec}_l(\mathbf{Y}, \mathbf{Z})$ with $\mathsf{Enc}_l(\mathbf{Z}, k)$ and $\mathsf{Dec}_l(\mathbf{Y}, \mathbf{Z}, k')$.

## C. Any-to-Many Conversion

With both the one-to-one and many-to-many VTNs, the source speaker must be known and specified at both training and test time. However, in some cases we would need to handle *any*-to-many VC tasks, namely, to convert the voice of an arbitrary speaker or an arbitrary speaking style that is not included in the training dataset. Any-to-many VC is attractive in that it allows for input speech of unknown speakers without the need for extra processing, such as model retraining and

adaptation. It may also be useful as a voice normalization preprocessor for speaker-independent ASR. Another important advantage of the many-to-many extension presented above is that it can be modified to handle any-to-many VC tasks by intentionally not allowing the source prenet and encoder to take the source speaker index $k$ as inputs. Namely, with this modified version, the output sequence of each layer in these networks is directly passed to the next layer without going through Eqs. (22) and (23). We show later how well this modified version performs on an any-to-many VC task in which the source speaker is unseen in the training dataset.

## D. Real-Time System Settings

It is important to be aware of real-time requirements when building VC systems. To let the VTN work in real-time, we need to make two modifications. First, the source prenet and encoder must not use future information, as with the target prenet, decoder, and postnet during training. This can be implemented by constraining the convolution layers in the source prenet to be causal and replacing Eq. (3) with Eq. (16) for all the sub-layers in the encoder. Second, since the speaking rate and rhythm of input speech cannot be changed at test time, all the TSA matrices are simply set to identity matrices so that the speaking rate and rhythm will be kept unchanged.

## V. EXPERIMENTS

### A. Experimental Settings

To confirm the effectiveness of the ideas proposed in Section IV, we conducted objective and subjective evaluation experiments. We used the CMU Arctic database [68], which consists of recordings of four speakers, clb (female), bdl (male), slt (female), and rms (male), reading the same 1,132 phonetically balanced English utterances. We used all these speakers for training and evaluation. Hence, there were a total of twelve different combinations of source and target speakers. For each speaker, the first 1,000 and last 32 sentences of the 1,132 sentences were used for training and evaluation, respectively. All the speech signals were sampled at 16 kHz. As already detailed in Subsection III-A, for each utterance, the spectral envelope, log $F_0$, coded aperiodicity, and voiced/unvoiced information were extracted every 8 ms using the WORLD analyzer [66]. 28 mel-cepstral coefficients (MCCs) were then extracted from each spectral envelope using the Speech Processing Toolkit (SPTK)[1]. The reduction factor $r$ was set to 3. Thus, $D = (28 + 3) \times 3 = 93$.

### B. Network Architecture Details

Dropouts with rate 0.1 were applied to the input sequences before being fed into the source and target prenets and the postnet only at training time. For the nonlinear activation function $\phi$ in each FFN sub-layer, we chose to use the GLU function since it yielded slightly better performance than the ReLU function. The two prenets and the postnet were each designed using three 1D dilated convolution layers with kernel size 5, each followed by a GLU activation function, where

---

[1]https://github.com/r9y9/pysptk

weight normalization [69] was applied to each layer. The channel number $d$ was set at 256 for the one-to-one VTN and 512 for the many-to-many VTN, respectively. The middle channel number $d'$ of each FFN sub-layer was set at 512 for the one-to-one VTN and 1024 for the many-to-many VTN, respectively.

### C. Hyperparameter Settings

$\lambda_{\mathsf{dal}}$ and $\lambda_{\mathsf{iml}}$ were set at 2000 and 1, respectively. $\nu$ was set at 0.3 for both the one-to-one and many-to-many VTNs. The $L_1$ norm $\|\mathbf{X}\|_1$ used in Eqs. (19) and (25) were defined as a weighted norm $\|\mathbf{X}\|_1 = \sum_{n=1}^{N} \frac{1}{r} \sum_{j=1}^{r} \sum_{i=1}^{31} \gamma_i |x_{ij,n}|$, where $x_{1j,n}, \ldots, x_{28j,n}$, $x_{29j,n}$, $x_{30j,n}$ and $x_{31j,n}$ denote the entries of $\mathbf{X}$ corresponding to the 28 MCCs, log $F_0$, coded aperiodicity and voiced/unvoiced indicator at time $n$, and the weights were set at $\gamma_1 = \cdots = \gamma_{28} = \frac{1}{28}$, $\gamma_{29} = \frac{1}{10}$, $\gamma_{30} = \gamma_{31} = \frac{1}{50}$, respectively. $L$ and $H$ were set to 4 and 4 for the many-to-many version, and 6 and 1 for the one-to-one version, respectively. All the networks were trained simultaneously with random initialization. Adam optimization [70] was used for model training where the mini-batch size was 16 and 30,000 iterations were run. We configured each mini-batch to consist of only utterances of the same source-target speaker pair. The learning rate and exponential decay rate for the first moment for Adam were set to $1.0 \times 10^{-4}$ and 0.9 for the many-to-many version with the PreLN architecture and to $5.0 \times 10^{-5}$ and 0.9 otherwise. All these hyperparameters were tuned on a subset of the ATR Japanese speech database [71], which consisted of 503 phonetically balanced sentences uttered by two male and two female speakers.

### D. Objective Performance Measures

The test dataset consisted of speech samples of each speaker reading the same sentences. Thus, the quality of a converted feature sequence could be assessed by comparing it with the feature sequence of the reference utterance.

*1) Mel-Cepstral Distortion:* Given two mel-cepstra, we can use the mel-cepstral distortion (MCD) to measure their difference. We used the average of the MCDs taken along the DTW path between converted and reference feature sequences as the objective performance measure for each test utterance. Note that a smaller MCD indicates better performance.

*2) Log $F_0$ Correlation Coefficient:* To evaluate the log $F_0$ contour of converted speech, we used the correlation coefficient between the predicted and target log $F_0$ contours [72] as the objective performance measure. In the experiment, we used the average of the correlation coefficients taken over all the test utterances as the objective performance measure for log $F_0$ prediction. Thus, the closer it is to 1, the better the performance. We call this measure the log $F_0$ correlation coefficient (LFC).

*3) Local Duration Ratio:* To evaluate the speaking rate and rhythm of converted speech, we used the measure called the local duration ratio (LDR) [47]. LDRs are defined as the local slopes of the DTW path between converted and reference utterances. In the following, we use the mean absolute difference between the LDRs and 1 (in percentage) as the overall
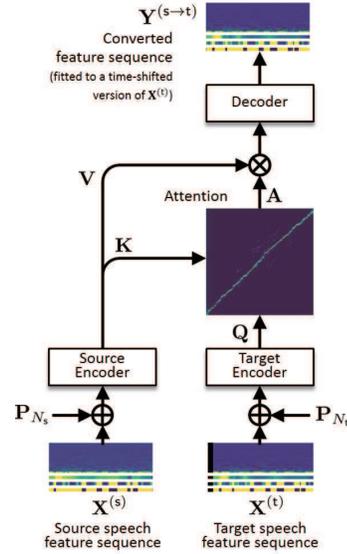


Fig. 8. Architecture of ConvS2S-VC

measure for the LDRs. Thus, the closer it is to zero, the better the performance. For example, if the converted speech is 2 times faster than the reference speech, the LDR will be 0.5 everywhere, and so its mean absolute difference from 1 will be $50\%$.

### E. Baseline Methods

*1) sprocket:* We chose the open-source VC method called sprocket [73] for comparison in our experiments. To run this method, we used the source code provided by the author[2]. Note that this method was used as a baseline in the Voice Conversion Challenge (VCC) 2018 [74].

*2) RNN-S2S-VC and ConvS2S-VC:* To compare different types of network architectures, we also tested the RNN-based S2S model [46], inspired by the architecture introduced in an S2S model-based TTS system called Tacotron [33], and the CNN-based model we presented previously [47]. We refer to these models as RNN-S2S-VC and ConvS2S-VC, respectively. **RNN-S2S-VC:** Although the original Tacotron used mel-spectra as the acoustic features, the baseline method was designed to use the same acoustic features as our method. The architecture was specifically designed as follows. The encoder consisted of a bottleneck fully connected prenet followed by a stack of $1 \times 1$ 1D GLU convolutions and a bi-directional LSTM layer. The decoder was an autoregressive content-based attention network consisting of a bottleneck fully connected prenet followed by a stateful LSTM layer producing the attention query, which was then passed to a stack of 2 uni-directional residual LSTM layers, followed by a linear projection to generate the features. **ConvS2S-VC:** The ConvS2S model we implemented for this experiment consisted of source/target encoders and a decoder, each of which had eight 1D GLU dilated convolution layers with kernel size of 5. We used single-step single-head scaled dot-product attention to compute attention distributions from

---

[2]https://github.com/k2kobayashi/sprocket

| Versions | | WA | Measures | | |
|---|---|---|---|---|---|
| | | | MCD(dB) | LFC | LDR(%) |
| one-to-one | PostLN | − | 7.12 | 0.705 | 5.75 |
| | | ✓ | 6.96 | 0.734 | 5.45 |
| | PreLN | − | 6.82 | 0.678 | 4.48 |
| | | ✓ | 6.66 | 0.703 | 3.90 |
| many-to-many | PostLN | − | 6.47 | 0.751 | 3.99 |
| | | ✓ | 6.35 | 0.761 | 3.90 |
| | PreLN | − | 6.28 | 0.759 | 4.16 |
| | | ✓ | 6.29 | 0.777 | 3.62 |

| Versions | | IML | Measures | | |
|---|---|---|---|---|---|
| | | | MCD(dB) | LFC | LDR(%) |
| many-to-many | PostLN | − | 6.94 | 0.644 | 4.12 |
| | | ✓ | 6.35 | 0.761 | 3.90 |
| | PreLN | − | 6.57 | 0.650 | 4.12 |
| | | ✓ | 6.28 | 0.792 | 2.51 |

the outputs of the source/target encoders. The convolutions in the target encoder and decoder were constrained to be causal, as with the target prenet and postnet in the one-to-one and many-to-many VTNs. A residual connection and weight normalization were applied to each layer in the three networks.

We designed and implemented many-to-many extensions of the RNN-based and CNN-based models to compare them with the many-to-many VTN.

### F. Objective Evaluations

*1) Ablation Studies:* We conducted ablation studies to confirm the effectiveness of the many-to-many VTN, IML, and attention windowing (AW), and compared the performances obtained with the PostLN and PreLN architectures. It should be noted that the models trained without the DAL were unsuccessful in producing recognizable speech, possibly due to the limited amount of training data. For this reason, we omit the results obtained when $\lambda_{dal} = 0$.

Tab. I lists the average MCDs, LFCs, and LDRs over the test samples obtained with the one-to-one and many-to-many VTNs with the PostLN and PreLN architectures with and without AW. Note that all the results for the many-to-many version were obtained with the models trained using the IML. We observed that the effect of the many-to-many VTN was noticeable. Comparisons between with and without AW revealed that while the AW showed a certain effect in improving the one-to-one version in terms of all the measures, it was found to be only slightly effective for the many-to-many version. This may imply that the prediction of attentions by the many-to-many version was already so successful that no correction by AW was necessary. As for the PostLN and PreLN architectures, the latter performed consistently better than the former, especially for the one-to-one version.

Tab. II shows the average MCDs, LFCs and LDRs over the test samples obtained with the many-to-many VTN trained with and without the IML. Note that all the results listed here

are obtained using AW. As these results indicate, the IML had a significant effect on performance improvements in terms of the MCD and LFC measures.

*2) Comparisons with Baseline Methods:* Tables III, IV and V show the MCDs, LFCs and LDRs obtained with the proposed and baseline methods. It should be noted that sprocket was designed to only adjust the mean and variance of the log $F_0$ contour of input speech and keep the rhythm unchanged. Hence, the performance gains over sprocket in terms of the LFC and LDR measures show how well the competing methods are able to predict the $F_0$ contours and rhythms of target speech. As the results indicate, all the S2S models performed better than sprocket in terms of the LFC and LDR measures, thus demonstrating the ability to properly convert the prosodic features in speech. They also performed better than or comparably to sprocket in terms of the MCD measure. It is worth noting that the many-to-many extension was found to be significantly effective for all the architecture types of the S2S models. It is interesting to compare the performance of the many-to-many versions of RNN-S2S-VC, ConvS2S-VC and VTN. The many-to-many version of ConvS2S-VC performed best in terms of the MCD and LFC measures whereas the many-to-many RNN-S2S-VC performed best in terms of the LDR measure. This may indicate that the strengths of S2S models can vary depending on the type of architecture.

As mentioned earlier, one important advantage of the transformer architecture over its RNN counterpart is that it can be trained efficiently thanks to its parallelizable structure. In fact, while it took about 30 hours and 50 hours to train the one-to-one and many-to-many versions of the RNN-S2S-VC model, it only took about 3 hours and 5 hours to train these two versions of the VTN under the current experimental settings. We implemented all the methods in PyTorch and used a single Tesla V100 GPU with a 32.0GB memory for training each model.

*3) Performance of any-to-many VTN:* We evaluated the performance of the any-to-many VTN described in Subsection IV-C under an open-set condition where the speaker of the test utterances is unseen in the training data. We used the utterances of speaker lnh (female) as the test input speech. The results are listed in Tab. VI (a). For comparison, Tab. VI (b) lists the results of sprocket performed on the same speaker pairs under a speaker-dependent closed-set condition. As these results indicate, the any-to-many VTN still performed better than sprocket, even though sprocket had an advantage in both the training and test conditions.

*4) Performance with Real-Time System Setting:* We evaluated the MCDs and LFCs obtained with the many-to-many VTN with the real-time system setting described in Subsection IV-D. The results are shown in Table VII. It is worth noting that it performed only slightly worse than with the default setting despite the restrictions related to the real-time system setting and performed still better than sprocket in terms of the MCD and LFC measures.

*5) Impact of training data size:* To evaluate the impact of the training data size, we compared the performance of the one-to-one and many-to-many VTNs trained using 1,000,

TABLE III
MCDs (dB) OBTAINED WITH BASELINE AND PROPOSED METHODS

| Speakers | | sprocket | RNN-S2S | | ConvS2S | | VTN | |
| source | target | | one-to-one | many-to-many | one-to-one | many-to-many | one-to-one | many-to-many |
|---|---|---|---|---|---|---|---|---|
| clb | bdl | 6.98 | 6.80 | 6.94 | 6.99 | 6.40 | 6.83 | 6.64 |
| | slt | 6.34 | 6.28 | 6.24 | 6.48 | 5.74 | 6.21 | 5.97 |
| | rms | 6.84 | 6.41 | 6.33 | 6.47 | 5.88 | 6.49 | 6.23 |
| bdl | clb | 6.44 | 6.33 | 6.20 | 6.61 | 5.79 | 6.17 | 6.03 |
| | slt | 6.46 | 6.49 | 6.24 | 6.68 | 5.92 | 6.45 | 6.13 |
| | rms | 7.24 | 6.53 | 6.28 | 6.76 | 6.09 | 6.80 | 6.34 |
| slt | clb | 6.21 | 6.20 | 6.21 | 6.41 | 5.69 | 6.20 | 5.91 |
| | bdl | 6.80 | 7.06 | 7.18 | 7.16 | 6.33 | 7.20 | 6.77 |
| | rms | 6.87 | 6.40 | 6.44 | 6.76 | 5.97 | 6.73 | 6.26 |
| rms | clb | 6.43 | 6.36 | 6.26 | 6.38 | 5.88 | 6.63 | 5.94 |
| | bdl | 7.40 | 7.07 | 7.13 | 7.40 | 6.56 | 7.51 | 6.74 |
| | slt | 6.76 | 6.47 | 6.29 | 6.71 | 6.01 | 6.75 | 6.21 |
| All pairs | | 6.73 | 6.50 | 6.39 | 6.74 | 6.03 | 6.66 | 6.29 |

TABLE IV
LFCs OBTAINED WITH BASELINE AND PROPOSED METHODS

| Speakers | | sprocket | RNN-S2S | | ConvS2S | | VTN | |
| source | target | | one-to-one | many-to-many | one-to-one | many-to-many | one-to-one | many-to-many |
|---|---|---|---|---|---|---|---|---|
| clb | bdl | 0.643 | 0.822 | 0.851 | 0.792 | 0.848 | 0.791 | 0.790 |
| | slt | 0.790 | 0.811 | 0.837 | 0.846 | 0.891 | 0.787 | 0.835 |
| | rms | 0.556 | 0.749 | 0.796 | 0.719 | 0.808 | 0.643 | 0.747 |
| bdl | clb | 0.642 | 0.738 | 0.809 | 0.752 | 0.828 | 0.785 | 0.767 |
| | slt | 0.632 | 0.768 | 0.837 | 0.808 | 0.871 | 0.703 | 0.784 |
| | rms | 0.467 | 0.716 | 0.716 | 0.732 | 0.801 | 0.595 | 0.759 |
| slt | clb | 0.820 | 0.748 | 0.774 | 0.795 | 0.849 | 0.750 | 0.772 |
| | bdl | 0.663 | 0.766 | 0.813 | 0.729 | 0.835 | 0.704 | 0.775 |
| | rms | 0.611 | 0.713 | 0.786 | 0.712 | 0.747 | 0.660 | 0.730 |
| rms | clb | 0.632 | 0.785 | 0.815 | 0.809 | 0.829 | 0.632 | 0.717 |
| | bdl | 0.648 | 0.816 | 0.833 | 0.799 | 0.814 | 0.689 | 0.810 |
| | slt | 0.674 | 0.777 | 0.805 | 0.811 | 0.804 | 0.715 | 0.748 |
| All pairs | | 0.653 | 0.775 | 0.808 | 0.785 | 0.826 | 0.703 | 0.777 |

TABLE V
LDR DEVIATIONS (%) OBTAINED WITH BASELINE AND PROPOSED METHODS

| Speakers | | sprocket | RNN-S2S | | ConvS2S | | VTN | |
| source | target | | one-to-one | many-to-many | one-to-one | many-to-many | one-to-one | many-to-many |
|---|---|---|---|---|---|---|---|---|
| clb | bdl | 17.66 | 3.52 | 3.22 | 3.04 | 3.34 | 3.33 | 2.34 |
| | slt | 9.74 | 2.34 | 2.70 | 0.86 | 4.38 | 3.90 | 4.28 |
| | rms | 3.24 | 2.70 | 3.76 | 4.18 | 6.79 | 5.57 | 2.82 |
| bdl | clb | 16.65 | 3.05 | 3.53 | 4.55 | 2.46 | 2.66 | 4.27 |
| | slt | 4.58 | 4.18 | 4.21 | 6.10 | 3.52 | 3.77 | 3.04 |
| | rms | 15.20 | 2.90 | 2.21 | 3.78 | 3.42 | 2.69 | 3.00 |
| slt | clb | 9.25 | 2.89 | 3.41 | 4.23 | 2.32 | 3.74 | 3.29 |
| | bdl | 5.52 | 2.35 | 2.04 | 3.47 | 3.72 | 3.69 | 3.03 |
| | rms | 11.46 | 3.06 | 5.00 | 3.66 | 5.57 | 6.88 | 5.41 |
| rms | clb | 2.84 | 4.50 | 4.17 | 3.91 | 1.54 | 4.28 | 3.35 |
| | bdl | 17.76 | 4.68 | 3.19 | 3.41 | 4.15 | 3.31 | 3.40 |
| | slt | 11.95 | 4.74 | 3.74 | 3.61 | 4.48 | 4.28 | 4.27 |
| All pairs | | 10.60 | 3.30 | 3.38 | 3.77 | 3.69 | 3.90 | 3.62 |

500, and 250 utterances for each speaker, respectively. These results are shown in Tables VIII, IX, and X. We can confirm that the performance of both versions degrades as expected in terms of all the measures as the training data size decreases. More importantly, we can see that the many-to-many VTN trained using only 500 utterances for each speaker performed comparably or slightly better than the one-to-one VTN and sprocket trained using 1,000 utterances for each speaker. This implies the fact that in training the mapping between a certain speaker pair using $500 \times 2$ utterances, $500 \times 2$ utterances of the remaining two speakers are as valuable as another $500 \times 2$ utterances of that speaker pair, if leveraged efficiently.

### G. Subjective Listening Tests

We conducted subjective listening tests to compare the sound quality and speaker similarity of the converted speech samples obtained with the proposed and baseline methods. For these tests, we used 32 speech samples generated by each method for each source-target speaker pair.

With the sound quality test, we evaluated the mean opinion score (MOS) for each speech sample. In this test, we included the speech samples synthesized in the same manner
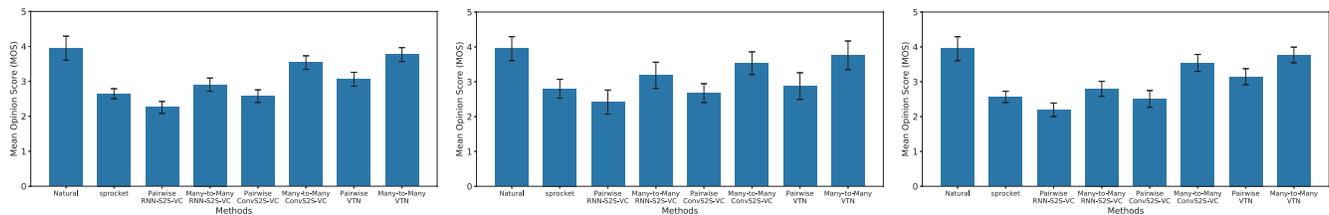
Fig. 9. Sound quality scores averaged across all speaker, intra-gender, and inter-gender pairs, respectively (from left to right).
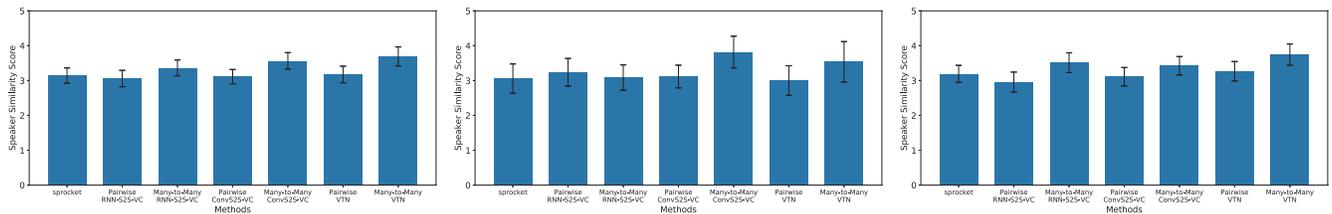


Fig. 10. Speaker similarity scores averaged across all speaker, intra-gender, and inter-gender pairs, respectively (from left to right).

TABLE VI

PERFORMANCE OF MANY-TO-MANY VTN WITH ANY-TO-MANY SETTING UNDER OPEN-SET CONDITION AND SPROCKET UNDER CLOSED-SET CONDITION TESTED ON SAME SAMPLES

(a) any-to-many VTN

| Speaker pair | | Measures | | |
|---|---|---|---|---|
| source | target | $MCD_{(dB)}$ | LFC | $LDR_{(\%)}$ |
| lnh | clb | 6.49 | 0.690 | 2.18 |
| | bdl | 7.24 | 0.636 | 4.44 |
| | slt | 6.59 | 0.693 | 4.40 |
| | rms | 6.87 | 0.466 | 8.65 |
| All pairs | | 6.71 | 0.630 | 4.41 |

(b) sprocket

| Speaker pair | | Measures | | |
|---|---|---|---|---|
| source | target | $MCD_{(dB)}$ | LFC | $LDR_{(\%)}$ |
| lnh | clb | 6.76 | 0.716 | 6.61 |
| | bdl | 8.26 | 0.523 | 13.38 |
| | slt | 6.62 | 0.771 | 5.72 |
| | rms | 7.22 | 0.480 | 4.87 |
| All pairs | | 7.21 | 0.579 | 7.61 |

TABLE VII

PERFORMANCE OF MANY-TO-MANY VTN WITH REAL-TIME SYSTEM SETTING

| Speaker pair | | Measures | |
|---|---|---|---|
| source | target | $MCD_{(dB)}$ | LFC |
| clb | bdl | 7.27 | 0.735 |
| | slt | 6.13 | 0.791 |
| | rms | 6.75 | 0.693 |
| bdl | clb | 6.36 | 0.685 |
| | slt | 6.61 | 0.715 |
| | rms | 6.61 | 0.660 |
| slt | clb | 6.12 | 0.743 |
| | bdl | 7.10 | 0.673 |
| | rms | 6.55 | 0.609 |
| rms | clb | 6.06 | 0.737 |
| | bdl | 7.22 | 0.612 |
| | slt | 6.60 | 0.730 |
| All pairs | | 6.58 | 0.703 |

as the proposed and baseline methods (namely, the WORLD synthesizer) using the acoustic features directly extracted from real speech samples. We also included speech samples produced using the one-to-one and many-to-many versions of

RNN-S2S-VC, ConvS2S-VC, and VTN and sprocket in the stimuli. Twenty listeners participated in our listening tests. Each listener was asked to evaluate the naturalness by selecting 5: Excellent, 4: Good, 3: Fair, 2: Poor, or 1: Bad for each utterance. The scores averaged across all, intra-gender, and inter-gender pairs are shown in Fig. 9. The one-to-one VTN performed better than sprocket and the one-to-one versions of the other S2S models. We also confirmed that the many-to-many extension had a significant effect in improving the audio quality of all the S2S models. It is worth noting that the many-to-many VTN performed better than all the competing methods including the many-to-many version of ConvS2S-VC, even though the many-to-many version of ConvS2S-VC was found to outperform the many-to-many VTN in terms of the MCD and LFC measures through the objective evaluation experiments, as reported earlier. According to the two-sided Mann-Whitney test performed on the MOS scores of the many-to-many VTN and each of the remaining methods, the $p$-values for all the pairs except for the many-to-many VTN and many-to-many ConvS2S-VC pair were less than 0.05, indicating that the many-to-many VTN performed significantly better than all the competing methods except the many-to-many ConvS2S-VC in terms of sound quality.

With the speaker similarity test, each listener was given a converted speech sample and a real speech sample of the corresponding target speaker and was asked to evaluate how likely they are to be produced by the same speaker by selecting 5: Definitely, 4: Likely, 3: Fairly likely, 2: Not very likely, or 1: Unlikely. We used converted speech samples generated by the one-to-one and many-to-many versions of RNN-S2S-VC and ConvS2S-VC, and sprocket for comparison as with the sound quality test. The scores averaged across all, intra-gender, and inter-gender pairs are shown in Fig. 10. The many-to-many versions of ConvS2S-VC and VTN performed comparably and slightly better than all other methods. According to the two-sided Mann-Whitney test, the many-to-many VTN was found to perform significantly better than the one-to-one VTN, one-to-one ConvS2S-VC, one-to-one RNN-S2S-VC, and sprocket in terms of speaker similarity.

TABLE VIII
MCDs (dB) OBTAINED WITH ONE-TO-ONE AND MANY-TO-MANY VTNs UNDER DIFFERENT TRAINING DATA SIZE CONDITIONS

| Speakers | | 1000 training utterances | | 500 training utterances | | 250 training utterances | |
|---|---|---|---|---|---|---|---|
| source | target | one-to-one | many-to-many | one-to-one | many-to-many | one-to-one | many-to-many |
| | bdl | 6.83 | 6.64 | 7.04 | 7.03 | 7.49 | 7.55 |
| clb | slt | 6.21 | 5.97 | 6.44 | 6.27 | 6.79 | 6.79 |
| | rms | 6.49 | 6.23 | 6.72 | 6.56 | 7.22 | 7.02 |
| | clb | 6.17 | 6.03 | 6.53 | 6.41 | 7.37 | 6.80 |
| bdl | slt | 6.45 | 6.13 | 6.72 | 6.48 | 7.19 | 6.97 |
| | rms | 6.80 | 6.34 | 7.24 | 6.71 | 7.67 | 7.40 |
| | clb | 6.20 | 5.91 | 6.36 | 6.28 | 6.92 | 6.67 |
| slt | bdl | 7.20 | 6.77 | 7.21 | 7.24 | 7.76 | 7.82 |
| | rms | 6.73 | 6.26 | 6.94 | 6.73 | 7.67 | 7.24 |
| | clb | 6.63 | 5.94 | 6.86 | 6.40 | 7.36 | 6.97 |
| rms | bdl | 7.51 | 6.74 | 7.39 | 7.13 | 8.11 | 7.83 |
| | slt | 6.75 | 6.21 | 7.23 | 6.45 | 7.65 | 7.04 |
| All pairs | | 6.66 | 6.29 | 6.89 | 6.64 | 7.47 | 7.06 |

TABLE IX
LFCs OBTAINED WITH ONE-TO-ONE AND MANY-TO-MANY VTNs UNDER DIFFERENT TRAINING DATA SIZE CONDITIONS

| Speakers | | 1000 training utterances | | 500 training utterances | | 250 training utterances | |
|---|---|---|---|---|---|---|---|
| source | target | one-to-one | many-to-many | one-to-one | many-to-many | one-to-one | many-to-many |
| | bdl | 0.791 | 0.790 | 0.758 | 0.689 | 0.709 | 0.676 |
| clb | slt | 0.787 | 0.835 | 0.801 | 0.827 | 0.815 | 0.803 |
| | rms | 0.643 | 0.747 | 0.693 | 0.675 | 0.631 | 0.574 |
| | clb | 0.785 | 0.767 | 0.797 | 0.725 | 0.721 | 0.719 |
| bdl | slt | 0.703 | 0.784 | 0.756 | 0.818 | 0.759 | 0.769 |
| | rms | 0.595 | 0.759 | 0.647 | 0.724 | 0.573 | 0.503 |
| | clb | 0.750 | 0.772 | 0.804 | 0.771 | 0.746 | 0.720 |
| slt | bdl | 0.704 | 0.775 | 0.715 | 0.665 | 0.578 | 0.639 |
| | rms | 0.660 | 0.730 | 0.658 | 0.655 | 0.567 | 0.533 |
| | clb | 0.632 | 0.717 | 0.682 | 0.704 | 0.591 | 0.673 |
| rms | bdl | 0.689 | 0.810 | 0.754 | 0.628 | 0.637 | 0.513 |
| | slt | 0.715 | 0.748 | 0.630 | 0.765 | 0.666 | 0.73 |
| All pairs | | 0.703 | 0.777 | 0.717 | 0.723 | 0.668 | 0.688 |

TABLE X
LDR DEVIATIONS (%) OBTAINED WITH ONE-TO-ONE AND MANY-TO-MANY VTNs UNDER DIFFERENT TRAINING DATA SIZE CONDITIONS

| Speakers | | 1000 training utterances | | 500 training utterances | | 250 training utterances | |
|---|---|---|---|---|---|---|---|
| source | target | one-to-one | many-to-many | one-to-one | many-to-many | one-to-one | many-to-many |
| | bdl | 3.33 | 2.34 | 12.10 | 8.47 | 9.03 | 12.30 |
| clb | slt | 3.90 | 4.28 | 4.87 | 2.66 | 4.71 | 4.16 |
| | rms | 5.57 | 2.82 | 3.22 | 3.73 | 4.24 | 7.66 |
| | clb | 2.66 | 4.27 | 7.25 | 3.67 | 4.69 | 2.61 |
| bdl | slt | 3.77 | 3.04 | 8.19 | 3.34 | 8.29 | 3.06 |
| | rms | 2.69 | 3.00 | 9.44 | 5.61 | 8.12 | 1.80 |
| | clb | 3.74 | 3.29 | 3.05 | 3.01 | 3.37 | 3.36 |
| slt | bdl | 3.69 | 3.03 | 10.15 | 8.49 | 11.51 | 9.73 |
| | rms | 6.88 | 5.41 | 5.18 | 2.26 | 7.59 | 3.44 |
| | clb | 4.28 | 3.35 | 2.94 | 3.21 | 5.33 | 4.47 |
| rms | bdl | 3.31 | 3.40 | 10.10 | 10.65 | 11.99 | 13.88 |
| | slt | 4.28 | 4.27 | 6.94 | 3.41 | 5.57 | 3.82 |
| All pairs | | 3.90 | 3.62 | 6.35 | 4.13 | 6.70 | 4.63 |

Audio samples of the one-to-one and many-to-many VTNs are available on the web[3].

## VI. CONCLUSION

We proposed an extension of VTN, which provides the flexibility of handling many-to-many, any-to-many, and real-time VC tasks without relying on ASR models and text annotations. Through ablation studies, we confirmed the effectiveness of each of the proposed ideas. Objective and subjective evaluation

experiments on a speaker identity conversion task showed that the proposed method could perform better than baseline methods.

Although we used the WORLD vocoder for waveform generation in the above experiments, using a neural vocoder instead could significantly improve the quality of the converted speech. Rather than simply performing feature mapping and then using a neural vocoder to generate waveforms, we believe that further improvements could be made by integrating the VTN and a neural vocoder into a single model so that the whole model can be trained end-to-end.

[3]http://www.kecl.ntt.co.jp/people/kameoka.hirokazu/Demos/vtn/index.html

Zero-shot VC is a task of converting input speech to the voice or speaking style of an unseen speaker by looking at only a few of his/her utterances [75]. Although in the many-to-many VTN, the target voice or speaking style is specified via a target speaker embedding vector, the embedding vector currently used for target speaker conditioning is nongeneralizable to unseen speakers. However, as proposed in [75], replacing the embedding vector with one used for speaker verification [76] may allow our model to handle zero-shot VC.

## REFERENCES

[1] A. Kain and M. W. Macon, "Spectral voice conversion for text-to-speech synthesis," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1998, pp. 285–288.

[2] A. B. Kain, J.-P. Hosom, X. Niu, J. P. van Santen, M. Fried-Oken, and J. Staehely, "Improving the intelligibility of dysarthric speech," *Speech Communication*, vol. 49, no. 9, pp. 743–759, 2007.

[3] K. Nakamura, T. Toda, H. Saruwatari, and K. Shikano, "Speaking-aid systems using GMM-based voice conversion for electrolaryngeal speech," *Speech Communication*, vol. 54, no. 1, pp. 134–146, 2012.

[4] Z. Inanoglu and S. Young, "Data-driven emotion conversion in spoken English," *Speech Communication*, vol. 51, no. 3, pp. 268–283, 2009.

[5] T. Toda, M. Nakagiri, and K. Shikano, "Statistical voice conversion techniques for body-conducted unvoiced speech enhancement," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 9, pp. 2505–2517, 2012.

[6] D. Felps, H. Bortfeld, and R. Gutierrez-Osuna, "Foreign accent conversion in computer assisted pronunciation training," *Speech Communication*, vol. 51, no. 10, pp. 920–932, 2009.

[7] Y. Stylianou, O. Cappé, and E. Moulines, "Continuous probabilistic transform for voice conversion," *IEEE Trans. SAP*, vol. 6, no. 2, pp. 131–142, 1998.

[8] T. Toda, A. W. Black, and K. Tokuda, "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 8, pp. 2222–2235, 2007.

[9] E. Helander, T. Virtanen, J. Nurminen, and M. Gabbouj, "Voice conversion using partial least squares regression," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 5, pp. 912–921, 2010.

[10] E. Helander, H. Silen, T. Virtanen, and M. Gabbouj, "Voice conversion using dynamic kernel partial least squares regression," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 3, pp. 806–817, 2011.

[11] X. Tian, S. W. Lee, Z. Wu, E. S. Chng, and H. Li, "An exemplar-based approach to frequency warping for voice conversion," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 10, pp. 1863–1876, 2017.

[12] R. Takashima, T. Takiguchi, and Y. Ariki, "Exemplar-based voice conversion in noisy environment," in *Proc. IEEE Spoken Language Technology Workshop (SLT)*, 2012, pp. 313–317.

[13] B. Sisman, H. Li, and K. C. Tan, "Sparse representation of phonetic features for voice conversion with and without parallel data," in *Proc. IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2017, pp. 677–684.

[14] S. Desai, A. W. Black, B. Yegnanarayana, and K. Prahallad, "Spectral mapping using artificial neural networks for voice conversion," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 5, pp. 954–964, 2010.

[15] S. H. Mohammadi and A. Kain, "Voice conversion using deep neural networks with speaker-independent pre-training," in *Proc. IEEE Spoken Language Technology Workshop (SLT)*, 2014, pp. 19–23.

[16] L. Sun, S. Kang, K. Li, and H. Meng, "Voice conversion using deep bidirectional long short-term memory based recurrent neural networks," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 4869–4873.

[17] H. Ming, D. Huang, L. Xie, J. Wu, M. Dong, and H. Li, "Deep bidirectional LSTM modeling of timbre and prosody for emotional voice conversion," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2016, pp. 2453–2457.

[18] C.-C. Hsu, H.-T. Hwang, Y.-C. Wu, Y. Tsao, and H.-M. Wang, "Voice conversion from non-parallel corpora using variational auto-encoder," in *Proc. Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2016, pp. 1–6.

[19] ——, "Voice conversion from unaligned corpora using variational autoencoding Wasserstein generative adversarial networks," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2017, pp. 3364–3368.

[20] Y. Saito, Y. Ijima, K. Nishida, and S. Takamichi, "Non-parallel voice conversion using variational autoencoders conditioned by phonetic posteriorgrams and d-vectors," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5274–5278.

[21] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "ACVAE-VC: Non-parallel voice conversion with auxiliary classifier variational autoencoder," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 9, pp. 1432–1443, 2019.

[22] P. L. Tobing, Y.-C. Wu, T. Hayashi, K. Kobayashi, and T. Toda, "Non-parallel voice conversion with cyclic variational autoencoder," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2019, pp. 674–678.

[23] T. Kaneko and H. Kameoka, "CycleGAN-VC: Non-parallel voice conversion using cycle-consistent adversarial networks," in *Proc. European Signal Processing Conference (EUSIPCO)*, 2018, pp. 2100–2104.

[24] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "StarGAN-VC: Non-parallel many-to-many voice conversion using star generative adversarial networks," in *Proc. IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 266–273.

[25] J. Serrà, S. Pascual, and C. Segura, "Blow: A single-scale hyperconditioned flow for non-parallel raw-audio voice conversion," *arXiv:1906.00794 [cs.LG]*, Jun. 2019.

[26] B. Sisman, J. Yamagishi, S. King, and H. Li, "An overview of voice conversion and its challenges: From statistical modeling to deep learning," *arXiv:2008.03648 [eess.AS]*, 2020.

[27] G. Sanchez, H. Silen, J. Nurminen, and M. Gabbouj, "Hierarchical modeling of F0 contours for voice conversion," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2014, pp. 2318–2321.

[28] H. Ming, D. Huang, L. Xie, J. Wu, M. Dong, and H. Li, "Exemplar-based sparse representation of timbre and prosody for voice conversion," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 5175–5179.

[29] Z. Luo, J. Chen, T. Takiguchi, and Y. Ariki, "Emotional voice conversion with adaptive scales F0 based on wavelet transform using limited amount of emotional data," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2017, p. 3399.3403.

[30] B. Sisman, M. Zhang, and H. Li, "Group sparse representation with WaveNet vocoder adaptation for spectrum and prosody conversion," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 6, pp. 1085–1097, 2019.

[31] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Adv. Neural Information Processing Systems (NIPS)*, 2014, pp. 3104–3112.

[32] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *Adv. Neural Information Processing Systems (NIPS)*, 2015, pp. 577–585.

[33] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, "Tacotron: Towards end-to-end speech synthesis," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2017, pp. 4006–4010.

[34] S. O. Arık, M. Chrzanowski, A. Coates, G. Diamos, A. Gibiansky, Y. Kang, X. Li, J. Miller, A. Ng, J. Raiman, S. Sengupta, and M. Shoeybi, "Deep Voice: Real-time neural text-to-speech," in *Proc. International Conference on Machine Learning (ICML)*, 2017.

[35] S. O. Arık, G. Diamos, A. Gibiansky, J. Miller, K. Peng, W. Ping, J. Raiman, and Y. Zhou, "Deep Voice 2: Multi-speaker neural text-to-speech," in *Adv. Neural Information Processing Systems (NIPS)*, 2017.

[36] J. Sotelo, S. Mehri, K. Kumar, J. F. Santos, K. Kastner, A. Courville, and Y. Bengio, "Char2Wav: End-to-end speech synthesis," in *Proc. International Conference on Learning Representations (ICLR)*, 2017.

[37] H. Tachibana, K. Uenoyama, and S. Aihara, "Efficiently trainable text-to-speech system based on deep convolutional networks with guided attention," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 4784–4788.

[38] W. Ping, K. Peng, A. Gibiansky, S. O. Arık, A. Kannan, S. Narang, J. Raiman, and J. Miller, "Deep Voice 3: Scaling text-to-speech with convolutional sequence learning," in *Proc. International Conference on Learning Representations (ICLR)*, 2018.

[39] M.-T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *Proc. EMNLP*, 2015.

[40] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," *arXiv:1705.03122 [cs.CL]*, May 2017.

[41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Adv. Neural Information Processing Systems (NIPS)*, 2017.

[42] H. Miyoshi, Y. Saito, S. Takamichi, and H. Saruwatari, "Voice conversion using sequence-to-sequence learning of context posterior probabilities," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2017, pp. 1268–1272.

[43] J.-X. Zhang, Z.-H. Ling, L.-J. Liu, Y. Jiang, and L.-R. Dai, "Sequence-to-sequence acoustic modeling for voice conversion," *arXiv:1810.06865 [cs.SD]*, Oct. 2018.

[44] M. Zhang, X. Wang, F. Fang, H. Li, and J. Yamagishi, "Joint training framework for text-to-speech and voice conversion using multi-source Tacotron and WaveNet," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2019, pp. 1298–1302.

[45] F. Biadsy, R. J. Weiss, P. J. Moreno, D. Kanevsky, 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. Annual Conference of the International Speech Communication Association (Interspeech)*, 2019, pp. 4115–4119.

[46] K. Tanaka, H. Kameoka, T. Kaneko, and N. Hojo, "AttS2S-VC: Sequence-to-sequence voice conversion with attention and context preservation mechanisms," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 6805–6809.

[47] H. Kameoka, K. Tanaka, D. Kwaśny, T. Kaneko, and N. Hojo, "ConvS2S-VC: Fully convolutional sequence-to-sequence voice conversion," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1849–1863, 2020.

[48] 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. Annual Conference of the International Speech Communication Association (Interspeech)*, 2020.

[49] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "Non-parallel voice conversion with augmented classifier star generative adversarial networks," *arXiv:2008.12604 [eess.AS]*, 2020.

[50] Q. Wang, B. Li, T. Xiao, J. Zhu, C. Li, D. F. Wong, and L. S. Chao, "Learning deep transformer models for machine translation," in *Proc. Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019, pp. 1810–1822.

[51] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T.-Y. Liu, "On layer normalization in the transformer architecture," in *Proc. International Conference on Machine Learning (ICML)*, 2020, pp. 503–512.

[52] T. Fukada, K. Tokuda, T. Kobayashi, and S. Imai, "An adaptive algorithm for mel-cepstral analysis of speech," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1992, pp. 137–140.

[53] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "WaveNet: A generative model for raw audio," *arXiv:1609.03499 [cs.SD]*, Sep. 2016.

[54] A. Tamamori, T. Hayashi, K. Kobayashi, K. Takeda, and T. Toda, "Speaker-dependent WaveNet vocoder," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2017, pp. 1118–1122.

[55] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. van den Oord, S. Dieleman, and K. Kavukcuoglu, "Efficient neural audio synthesis," *arXiv:1802.08435 [cs.SD]*, Feb. 2018.

[56] S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. Courville, and Y. Bengio, "SampleRNN: An unconditional end-to-end neural audio generation model," *arXiv:1612.07837 [cs.SD]*, Dec. 2016.

[57] Z. Jin, A. Finkelstein, G. J. Mysore, and J. Lu, "FFTNet: A real-time speaker-dependent neural vocoder," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 2251–2255.

[58] A. van den Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. van den Driessche, E. Lockhart, L. C. Cobo, F. Stimberg, N. Casagrande, D. Grewe, S. Noury, S. Dieleman, E. Elsen, N. Kalchbrenner, H. Zen, A. Graves, H. King, T. Walters, D. Belov, and D. Hassabis, "Parallel WaveNet: Fast high-fidelity speech synthesis," *arXiv:1711.10433 [cs.LG]*, Nov. 2017.

[59] W. Ping, K. Peng, and J. Chen, "ClariNet: Parallel wave generation in end-to-end text-to-speech," *arXiv:1807.07281 [cs.CL]*, Feb. 2019.

[60] R. Prenger, R. Valle, and B. Catanzaro, "WaveGlow: A flow-based generative network for speech synthesis," *arXiv:1811.00002 [cs.SD]*, Oct. 2018.

[61] S. Kim, S. Lee, J. Song, and S. Yoon, "FloWaveNet: A generative flow for raw audio," *arXiv:1811.02155 [cs.SD]*, Nov. 2018.

[62] X. Wang, S. Takaki, and J. Yamagishi, "Neural source-filter-based waveform model for statistical parametric speech synthesis," *arXiv:1810.11946 [eess.AS]*, Oct. 2018.

[63] K. Tanaka, T. Kaneko, N. Hojo, and H. Kameoka, "Synthetic-to-natural speech waveform conversion using cycle-consistent adversarial networks," in *Proc. IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 632–639.

[64] K. Kumar, R. Kumar, T. de Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. de Brebisson, Y. Bengio, and A. Courville, "MelGAN: Generative adversarial networks for conditional waveform synthesis," in *Adv. Neural Information Processing Systems (NeurIPS)*, 2019, pp. 14 910–14 921.

[65] R. Yamamoto, E. Song, and J.-M. Kim, "Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6199–6203.

[66] 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. E99-D, no. 7, pp. 1877–1884, 2016.

[67] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, "Tacotron: Towards end-to-end speech synthesis," in *Proc. Annual Conference of the International Speech Communication Association (Interspeech)*, 2017, pp. 4006–4010.

[68] J. Kominek and A. W. Black, "The CMU Arctic speech databases," in *Proc. ISCA Speech Synthesis Workshop (SSW)*, 2004, pp. 223–224.

[69] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," in *Adv. Neural Information Processing Systems (NIPS)*, 2016, pp. 901–909.

[70] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. International Conference on Learning Representations (ICLR)*, 2015.

[71] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis," *Speech Communication*, vol. 9, no. 4, pp. 357–363, Aug. 1990.

[72] D. J. Hermes, "Measuring the perceptual similarity of pitch contours," *J. Speech Lang. Hear. Res.*, vol. 41, no. 1, pp. 73–82, 1998.

[73] K. Kobayashi and T. Toda, "sprocket: Open-source voice conversion software," in *Proc. Odyssey*, 2018, pp. 203–210.

[74] J. Lorenzo-Trueba, J. Yamagishi, T. Toda, D. Saito, F. Villavicencio, T. Kinnunen, and Z. Ling, "The voice conversion challenge 2018: Promoting development of parallel and nonparallel methods," *arXiv:1804.04262 [eess.AS]*, Apr. 2018.

[75] K. Qian, Y. Zhang, S. Chang, X. Yang, and M. Hasegawa-Johnson, "AutoVC: Zero-shot voice style transfer with only autoencoder loss," in *Proc. International Conference on Machine Learning (ICML)*, 2019, pp. 5210–5219.

[76] L. Wan, Q. Wang, A. Papir, and I. L. Moreno, "Generalized end-to-end loss for speaker verification," in *Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 4879–4883.