

TENSOR-TO-VECTOR REGRESSION FOR MULTI-CHANNEL SPEECH ENHANCEMENT BASED ON TENSOR-TRAIN NETWORK

Jun Qi^{1*}, Hu Hu^{1*}, Yannan Wang³, Chao-Han Huck Yang¹, Sabato Marco Siniscalchi^{1,2}, Chin-Hui Lee¹

¹Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA, USA

²Computer Engineering School, University of Enna, Italy

³Tencent Media Lab, Tencent Corporation, Shenzhen, Guangdong, China

ABSTRACT

We propose a tensor-to-vector regression approach to multi-channel speech enhancement in order to address the issue of input size explosion and hidden-layer size expansion. The key idea is to cast the conventional deep neural network (DNN) based vector-to-vector regression formulation under a tensor-train network (TTN) framework. TTN is a recently emerged solution for compact representation of deep models with fully connected hidden layers. Thus TTN maintains DNN's expressive power yet involves a much smaller amount of trainable parameters. Furthermore, TTN can handle a multi-dimensional tensor input by design, which exactly matches the desired setting in multi-channel speech enhancement. We first provide a theoretical extension from DNN to TTN based regression. Next, we show that TTN can attain speech enhancement quality comparable with that for DNN but with much fewer parameters, e.g., a reduction from 27 million to only 5 million parameters is observed in a single-channel scenario. TTN also improves PESQ over DNN from 2.86 to 2.96 by slightly increasing the number of trainable parameters. Finally, in 8-channel conditions, a PESQ of 3.12 is achieved using 20 million parameters for TTN, whereas a DNN with 68 million parameters can only attain a PESQ of 3.06. Code is available online¹.

Index Terms— Tensor-Train network, speech enhancement, deep neural network, tensor-to-vector regression

1. INTRODUCTION

Deep neural network (DNN) based speech enhancement [1] has demonstrated state-of-the-art performances in a single-channel setting. It has also been extended to multi-channel speech enhancement with similar high-quality enhanced speech [2]. A recent overview can be found in [3]. In essence, the process can be abstracted as a vector-to-vector regression based on deep architectures with a DNN aiming at learning a functional relationship $f: \mathbb{Y} \rightarrow \mathbb{X}$ such that input noisy speech $y \in \mathbb{Y}$ can be mapped to corresponding clean speech $x \in \mathbb{X}$. Several variants of deep learning structures have also been attempted, e.g., recurrent neural networks (RNNs) with long short term memory (LSTM) gates were employed in [4, 5]. A deep bidirectional RNN with LSTM gates was instead used in [6]. Moreover, a generative adversarial network (GAN) was used in [7].

Spatial information can complement spectral information for improved speech enhancement leveraging multi-channel information in more complex scenarios, where the sources of distortions include ambient noises, room reverberation and interfering speakers,

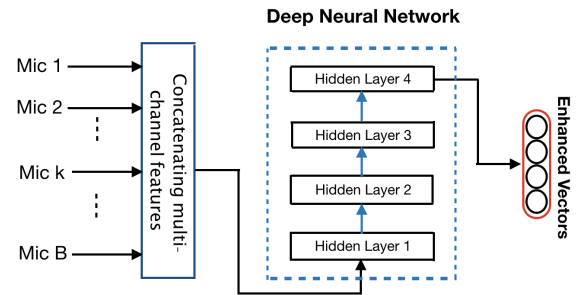


Fig. 1. Conventional multi-channel DNN-based vector-to-vector regression for speech enhancement.

as discussed in [8, 9, 10, 11], for example. However, the DNN-based vector-to-vector regression, which is the focus in this work, mainly aims at single-channel speech enhancement and is not simply generalized to multi-channel speech enhancement. As shown in Figure 1, a traditional approach to dealing with an array of microphones is exploited spatial information at the input level by concatenating speech vectors from multiple microphones into a single high dimensional vector, e.g., [12, 2]. Thus the vector-to-vector regression approach can still be employed for speech enhancement by appending multi-channel feature vectors together into a high-dimensional vector and mapping it to a vector extracted from the reference vector. Such a simple solution, unfortunately, clashes with our theoretical analysis outlined in [13], which suggests that the width of each hidden layer in a DNN needs to be greater than the input dimension plus two so that the expressive power of DNN-based vector-to-vector regression can be guaranteed. Moreover, high dimensional input vectors result in wider hidden non-linear layers, which in turn implies the need for huge computational resources and storage cost. Several proposals were put forth by different groups to overcome such an issue, for example, sparseness in DNN was explored in [14] to reduce the model size; however, the optimum implementation heavily depends on the specific hardware architecture. In [15], a singular value decomposition approach was devised, which is performed in a post-processing phase, that is, a DNN with a huge amount of parameters is first trained, and parameter reduction is then applied.

In this work, we proposed to address the multi-channel speech enhancement problem within a more principled parameter reduction framework, namely the Tensor-Train Network (TTN) [16], which does not require a multi-stage approach. A Tensor-Train refers to a compact tensor representation of the fully-connected hidden layers in a DNN [16]. In doing so, we put forth a tensor-to-vector

* Refers to an equal contribution.

¹<https://github.com/uwjunqi/Tensor-Train-Neural-Network>

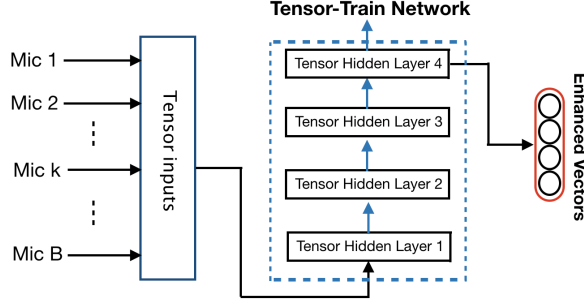


Fig. 2. A TTN-based speech enhancement system.

regression approach that can handle multi-channel information and meanwhile addresses the issue of input size explosion and hidden-layer size expansion. As shown in Figure 2, the feature vectors of B microphones are transformed into a tensor representation, such that the DNN-based vector-to-vector regression can be reshaped to a TTN based tensor-to-vector mapping. We show how the theorem on expressive power of the traditional DNN-based vector-to-vector [13] can be extended to be valid in the TTN framework. Next, we evaluate our approach with both single- and multi-microphone experimental setups, where the speech signal is exposed to multiple sources of distortions, including ambient noises, room reverberation and interfering speakers. We show that TTNs can attain speech enhancement quality comparable with DNNs, yet a 81% reduction in the size of parameters is observed, in the single-channel scenario. In 8-channel conditions, TTN delivers a PESQ of 3.12 using 20 million parameters. In contrast, a DNN-based multi-channel configuration needs up to 68 million parameters and attain a PESQ of 3.06. We should also remark that Tensor-Train decomposition can be applied to CNNs, and RNNs; however, we focus on DNNs in this initial investigation.

The remainder of the paper is organized as follows: Section 2 introduces TTN mathematical underpinnings and key properties. In Section 4, the experimental environment is first described, and then the experimental results are presented and discussed. Section 5 concludes our work.

2. TENSOR-TRAIN NETWORK

TTN relies on a tensor-train decomposition [17] which can be described mathematically as follows: given a vector of ranks $\mathbf{r} = \{r_1, r_2, \dots, r_{K+1}\}$, tensor-train decomposes a tensor $W \in \mathbb{R}^{(m_1 n_1) \times (m_2 n_2) \times \dots \times (m_K n_K)}$, $\forall i \in \{1, \dots, K\}$, $m_i \in \mathbb{R}^+$, $n_i \in \mathbb{R}^+$ into a multiplication of core tensors according to Eq. (1), where for the given ranks r_k and r_{k+1} , the k -th core tensor $C^{[k]}(r_k, i_k, j_k, r_{k+1}) \in \mathbb{R}^{m_k \times n_k}$ in which $i_k \in \{1, 2, \dots, m_k\}$ and $j_k \in \{1, 2, \dots, n_k\}$. Besides, r_1 and r_{K+1} are fixed to 1.

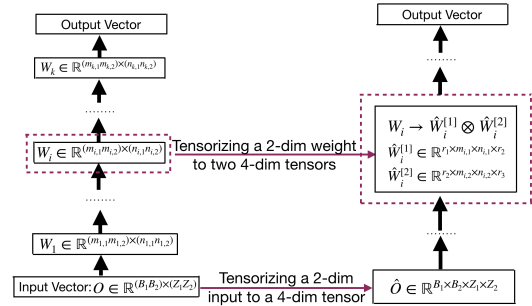
$$W((i_1, j_1), (i_2, j_2), \dots, (i_K, j_K)) = \prod_{k=1}^K C^{[k]}(r_k, i_k, j_k, r_{k+1}) \quad (1)$$

TTN is generated by applying the Tensor-Train decomposition to the hidden layers of DNN. The key benefit from TTN is to significantly reduce the number of parameters of a feed-forward DNN with fully-connected hidden layers. It is because TTN only stores the TT-format of DNN, i.e., the set of low-rank core tensors $\{C_k\}_{k=1}^K$

of the size $\sum_{k=1}^K m_k n_k r_k r_{k+1}$, which can approximately reconstruct the original DNN. In contrast, the memory storage of original DNN requires the size of $\prod_{k=1}^K m_k n_k$, which is much larger than $\sum_{k=1}^K m_k n_k r_k r_{k+1}$. Moreover, instead of decomposing a well-trained DNN into a TTN, core tensors can be randomly initialized and iteratively trained from scratch. Similar to the DNN training procedures, the optimization methods based on variants of stochastic gradient descent (SGD) can ensure converged solutions. Besides, the work [16] shows that the running complexities of DNN and the corresponding TTN are in the same order scale.

3. TENSOR-TO-VECTOR REGRESSION

The TTN framework offers a natural way to convert a speech enhancement system from vector-to-vector to tensor-to-vector configuration because multi-dimensional inputs can be directly fed into the TTN avoiding concatenating speech vectors inputs into a single long vector. Figure 3 demonstrates a typical example of casting DNN-based vector-to-vector regression into a TTN-based tensor-to-vector regression. More specifically, given a list of ranks $\{r_1, r_2, r_3\}$, the DNN input vector $O \in \mathbb{R}^{(B_1 B_2) \times (Z_1 Z_2)}$ is decomposed into a 4th order input tensors $\hat{O} \in \mathbb{R}^{B_1 \times B_2 \times Z_1 \times Z_2}$ before being used in the TTN-based configuration. Accordingly, the i -th weight matrix, W_i , in the DNN is decomposed into a tensor product of two tensors $\hat{W}_i^{[1]} \in \mathbb{R}^{r_1 \times m_{i,1} \times n_{i,1} \times r_2}$ and $\hat{W}_i^{[2]} \in \mathbb{R}^{r_2 \times m_{i,2} \times n_{i,2} \times r_3}$. The corresponding hidden layers with core tensors represent a basic architecture of TTN for tensor-to-vector regression, and the output vectors of two architectures are set to be equal.



(a) DNN for Vector-to-Vector Regression (b) TTN for Tensor-to-Vector Regression

Fig. 3. Transforming a DNN-based vector-to-vector regression into a TTN-based tensor-to-vector regression.

Tensor-to-vector regression based on TTN can substantially reduce the number of model parameters. Besides, we can prove that the TTN-based tensor-to-vector regression can maintain the representation power of the traditional DNN-based vector-to-vector. Our theoretical work [13] on the representation power of DNN-based vector-to-vector regression suggests a tight upper bound to a numerically estimate of the maximum mean squared estimation (MSE) loss. That MSE upper bound explicitly links the depth of ReLU-based hidden layers to the expressive capability of DNN models. More specifically, with input and out dimensions of d and q , respectively for a vector-to-vector regression target function $\hat{f}: \mathbb{R}^d \rightarrow \mathbb{R}^q$, there exists a DNN f_{DNN} with $K(K \geq 2)$ ReLU based hidden layers, where the width of each hidden layer is at least $(d+2)$ and the top hidden layer has $n_K(n_K \geq d+2)$ units. Then, we can derive the equality of (b) in Eq. (2).

$$\|\hat{f} - f_{TTN}\|_1 \stackrel{(a)}{\leq} \|\hat{f} - f_{DNN}\|_1 \stackrel{(b)}{=} O\left(\frac{q}{(n_K + K - 1)^2}\right) \quad (2)$$

We are currently developing new theories on TTN-based tensor-to-vector mapping, which demonstrate that a TTN associated with the tensor decomposition of a DNN can keep the expressive power of DNN-based vector-to-vector regression. Indeed, we can mathematically find a TTN f_{TTN} representation such that the inequality of (a) in Eq. (2) is valid. We avoid giving the proof of our new results here, but the experimental evidence given in the next section on the speech enhancement task supports our theoretical results in Eq. (2).

The inequality in Eq. (2) is ensured to be valid under consistent training and testing environmental conditions. It basically states that a TTN has more expressive power than a DNN with a similar complexity. The detailed proof is shown in our another paper [18], and the effect is demonstrated in the experimental results in Tables 1 and 2 to be provided later.

4. EXPERIMENTAL SETUP & RESULTS

4.1. Data Preparation

The proposed TTN-based models are evaluated on simulated data from WSJ0, which contains additive noise, interfering speakers, and reverberation. The dataset is created by corrupting the WSJ0 corpus [19] with OSU-100-noise [20] data, allowing us to obtain 30 hours of training material, and 5 hours of testing data. When simulating the noisy data, each waveform is mixed with one kind of background noise from the noise set. The target and additional interfering speech with their corresponding RIRs are convolved to generate the final waveform. In particular, our training and testing datasets are created from different noisy utterances of various speakers. As for the training dataset, a 5-minute clean speech from each of the targeted speakers is randomly mixed with 73 interfering speakers and 90 types of additive noise. Each target an-echoic speech is generated by convolving clean speech with the direct path response between the target speakers and the reference channel. For the testing dataset, another 5-minute unseen speech of the target speakers is mixed with 10 unseen interfering speakers and 10 types of unseen noise. The signal-to-interfered-noise-ratio (SINR) level of each utterance is set as follows: when SINR is 5dB, the signal-to-noise-ratio (SNR) is configured to 10dB or 15dB; when SINR is 10dB, the SNR is set to 15dB; when SINR is 15dB, the SNR is increased to 20dB. The proportion of each SINR level is set equally. Besides, some utterances of SINR 30dB are included in the training set to cover some very high SINR conditions.

To simulate reverberated speech, a reverberated acoustic environment is built: A microphone array of 8-circular channel microphone is arranged in a room of size $6.5m \times 5.5m \times 3m$ in terms of length-width-height. As to the single-channel scenario, the microphone is put at the center of the array. To avoid unnecessary combinations of multiple interference, we deliberately constrain the conditions that the microphone array only aims at one target speaker, and it only received one kind of additive noise. Specifically, a horizontal distance of a target speaker to the center of the microphone array is strictly fixed to $3m$. Besides, we set both the target speaker and the interfering speaker keeping the same distance to the microphone array, and the angle of them is configured as 40° . Before we build the training and testing sets, an improved image-source method (ISM) [21] is used to generate RIRs of reverberation time (RT60) (from $0.2s$ to $0.3s$) and the corresponding direct path response for

each microphone channel. For both training and testing datasets, the setting of RIRs is fixed to the same conditions, such as the room size, RT60, and all of the distances and directions.

4.2. Experiment Setting

In our experiments, 257-dimensional normalized log-power spectral (LPS) feature [22, 23, 24] is taken as the inputs to the DNNs. The LPS features are generated by computing 512 points Fourier transform on a speech segment of 32 milliseconds. For B -channel data, all channel inputs are concatenated together for the model training. For each input frame, the adjacent context of size M , is combined with the current frame. Hence, the size of the input for DNN is $257 \times (2M + 1) \times B$. For TTN, we ignore the first dimension of the input LPS features, which is the direct-current component. Thus, the input size for TTN is $256 \times (2M + 1) \times B$. After regression, the first dimension of input is concatenated back to the 256-dimension output without any change. The clean speech features of the first channel are assigned to the top layers of DNN and TTN, as the reference during the training stage.

Our baseline DNN model has 6 hidden layers, and each hidden layer is composed of 2048 neurons. The ReLU activation function is used for all hidden layers. A linear function is employed in the top hidden layer. As for the TTN, each hidden dense layer is decomposed and replaced with a Tensor-Train format, as shown in Figure 2. Both the DNN and TTN are trained from scratch based on the standard back-propagation algorithm [25], and both models adopt the same training configuration. During the training phase, Adam optimizer is adopted, and the initial learning rate is set to 0.0002. The mean square error (MSE) [26, 27] is used as the optimization object. The context window size at the input layer is set to 5 for all models, in which the current frame is concatenated with the previous 5 frames and following 5 frames (from the same channel). For both the DNN and TTN models, TensorFlow [28] is used in all of our experiments.

Perceptual evaluation of speech quality (PESQ) [29] is used as our evaluation criterion. The PESQ score, which ranges from -0.5 to 4.5 , is calculated by comparing the enhanced speech with the clean one. A higher PESQ score corresponds to a higher quality of the speech perception.

4.3. Single-channel Speech Enhancement Results

The TTN framework is first evaluated on the single-channel speech enhancement task. In Table 1, we can see that a 6-layer DNN model, which is taken as a baseline system, achieves a PESQ score of 2.86 with 27 million parameters. A TTN architecture is generated by applying tensor-train decomposition to the DNN baseline model. Each weight matrix of DNN is decomposed to two four-dimension tensors using tensor decomposition. For example, a weight matrix of the size 2048×2048 can be decomposed to two tensors with the size of $1 \times 32 \times 32 \times 4$ and $4 \times 64 \times 64 \times 1$. The TTN core tensors are randomly initialized and then trained from scratch using the Adam optimizer [30].

As discussed in Sections 2 and 3, the forth row in Table 1 shows that a substantial parameter reduction when the Tensor-Train format is employed. A drop in the PESQ value, from 2.86 (DNN) to 2.66 (TNN), is observed because of the parameter reduction. Nonetheless, TTNs consistently delivers better and better speech enhancement results as its number of parameters increases. As shown in the seventh row in Table 1, the TTN model with 5 million parameters can achieve nearly the same PESQ scores as the DNN baseline

model (2.84 vs 2.86). However, TTN model uses only 18% of the amount of parameters in the DNN, which suggests that TTN can significantly reduce the number of parameters while keeping the baseline performance. Furthermore, if we further increase the parameter number of the TTN model, we can even obtain a better TTN model achieving a 0.1 absolute PESQ improvement using only 74% of the DNN parameters, namely 20 million (see last row in the table).

Table 1. PESQ results for single-channel speech enhancement.

Model	Channel #	Parameter #	PESQ
DNN	1	27M	2.86
DNN-SVD	1	5M	2.82
DNN-SVD	1	20M	2.84
TTN	1	0.6M	2.66
TTN	1	0.7M	2.71
TTN	1	2.6M	2.78
TTN	1	5M	2.84
TTN	1	20M	2.96

To better appreciate the TTN technique, we have implemented the DNN-SVD method proposed in [15], which is a widely used post-processing compression technique for deep models. Each weight matrix in the trained DNN is first decomposed into two smaller matrix using SVD. The reduced size DNN is then fine-tuned with back-propagation. The DNN-SVD method can allow us to reduce the DNN size to 5 and 20 million parameters, as shown in the second and third rows in Table 1, respectively. Nonetheless, the PESQ drops to 2.82 with 5 million parameters. The PESQ attains a value of 2.84 when 20 million parameters are kept after DNN-SVD, which is however still lower than the PESQ attained by the original DNN. Most importantly, our TTN approach with 20 million parameters attains a PESQ of 2.96, which not only compares favourably with the DNN-SVD solution with the same amount of parameters, but it represents an improvement over the DNN baseline approach, as already mentioned.

4.4. Multi-channel Speech Enhancement Results

In this section, our investigation focuses on the microphone array scenario. The DNN-based multi-channel speech enhancement solutions realized in Figure 1 handles information from different mics at the input layer, which calls for wider input vectors and thereby an overall larger model size. By comparing the DNN baseline configuration in the first row in Tables 1 and 2, the parameter number goes from 27 million for the 1-channel case to 33 million for 2-channel case. In the 8-channel case, the number of DNN parameters goes up to 68 million, as shown in the third row in Table 2. In multi-channel speech enhancement, the control of the size of the neural architecture therefore becomes even a more important aspect that cannot be overlooked.

In the 2-channel configuration, TTN parameters can be reduced to 0.6 million, but the PESQ would go down to 2.75, as shown in the fourth row in Table 2. Nevertheless, we can observe, by comparing the first and sixth row in Table 2, that the TTN-based speech enhancement model can attain nearly the same performance of the DNN-based counterpart (3.00 vs. 2.96), while squeezing the number of parameters from 33 million to 5 million. As in the single-channel case, DNN-SVD could be applied as a post-processing approach for model compression, and it can achieve a PESQ of 2.93 with 5 million parameters, as shown in the fifth row in Table 2.

Table 2. PESQ Results of Multi-channel Speech Enhancement

Model	Channel #	Parameter #	PESQ
DNN	2	33M	3.00
DNN-SVD	2	5M	2.93
DNN	8	68M	3.06
DNN-SVD	8	5M	3.01
TTN	2	0.6M	2.75
TTN	2	5M	2.96
TTN	8	0.6M	2.83
TTN	8	5M	3.06
TTN	8	20M	3.12

In the 8-channel configuration, the TTN model size can be significantly reduced to 7% of the corresponding DNN size, namely from 68 million to 5 million, as shown in the eighth row in Table 2, while keeping the PESQ equal to 3.06. When compared with the result of DNN-SVD method [15] reported in the fourth row in Table 2, although both DNN-SVD and TTN can reduce the model parameters to 5 million, the PESQ will decrease from 3.06 to 3.01 with the DNN-SVD method; whereas the TTN approach keeps the model performance. Such an empirical result demonstrates the advantages in using a TTN-based tensor-to-vector regression approach to multi-channel speech enhancement and connects our theoretical analysis in Section 3. Furthermore, better speech enhancement results can be attained by increasing the TTN model size. Indeed, a TTN model with 20 million parameters, compared to the 68 million of the DNN baseline system, can achieve a PESQ of 3.12, as shown in the last row in Table 2, which corresponds to a 0.6 absolute improvement over the DNN baseline solution.

Since complex noisy and reverberated conditions are involved, the theory on the expressive power as shown in Eq. (2), is not strictly valid in our testing cases. But our improved theorems, which are to be published elsewhere and consider the generalization capability of models, can better correspond to our experimental results in Table 1 and 2.

5. CONCLUSIONS

This work is concerned with multi-channel speech enhancement leveraging upon TTN for tensor-to-vector regression. First, we describe the tensor-train decomposition, and its application to a DNN structure. Our theoretical results justify that the TTN-based tensor-to-vector function allows fewer parameters to realize a higher expressive and generalization power of DNN-based vector-to-vector regression. Then, we set up both single-channel and multi-channel speech enhancement systems based on TTN to verify our claims. When compared with DNN-based speech enhancement systems, experimental evidence shows that tensor-to-vector regression based on TTN compares favorably with the DNN-based baseline counterpart in terms of PESQ, but it requires much fewer parameters. Furthermore, increasing number in the TTN-based configuration, we can achieve better enhancement qualities in terms of PESQ, while amount of parameters used in the TTN configuration is much less than the DNN counterpart. The related empirical results verify our theories and exhibit a potential usage of tensor-train decomposition in more advanced deep learning models. Our future work will apply TTN decomposition to other deep models, such as recursive neural networks (RNNs) and convolution neural networks (CNNs).

6. REFERENCES

- [1] Y. Xu, J. Du, L.-R. Dai, and C.-H. Lee, "A regression approach to speech enhancement based on deep neural networks," *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)*, vol. 23, no. 1, pp. 7–19, 2015.
- [2] Q. Wang, S. Wang, F. Ge, C. W. Han, J. Lee, L. Guo, and C.-H. Lee, "Two-stage enhancement of noisy and reverberant microphone array speech for automatic speech recognition systems trained with only clean speech," in *ISCSLP*, 2018, pp. 21–25.
- [3] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 10, pp. 1702–1726, 2018.
- [4] F. Weninger, H. Erdogan, S. Watanabe, E. Vincent, J. Le Roux, J. R. Hershey, and B. Schuller, "Speech enhancement with lstm recurrent neural networks and its application to noise-robust asr," in *International Conference on Latent Variable Analysis and Signal Separation*. Springer, 2015, pp. 91–99.
- [5] H. Zhao, S. Zarar, I. Tashev, and C.-H. Lee, "Convolutional-recurrent neural networks for speech enhancement," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2018, pp. 2401–2405.
- [6] L. Sun, S. Kang, K. Li, and H. Meng, "Voice conversion using deep bidirectional long short-term memory based recurrent neural networks," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2015, pp. 4869–4873.
- [7] S. Pascual, A. Bonafonte, and J. Serra, "Segan: Speech enhancement generative adversarial network," *arXiv preprint arXiv:1703.09452*, 2017.
- [8] H. L. Van Trees, *Optimum array processing: Part IV of detection, estimation, and modulation theory*, John Wiley and Sons., 2004.
- [9] J. Benesty, J. Chen, and Y. Huang, *Microphone array signal processing*, JohSpringer Science and Business Media, 2008.
- [10] X. Xiao, S. Watanabe, H. Erdogan, L. Lu, J. Hershey, M. L. Seltzer G. Chen, Y. Zhang, M. Mandel, and D. Yu, "Deep beamforming networks for multi-channel speech recognition," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2016, pp. 5745–5749.
- [11] Z. Wang and D. Wang, "All-neural multi-channel speech enhancement," in *Interspeech*. ISCA, 2018, pp. 3234–3238.
- [12] B. Wu, M. Yang, K. Li, Z. Huang, S.M. Siniscalchi, T. Wang, and C.-H. Lee, "A reverberation-time-aware DNN approach leveraging spatial information for microphone array dereverberation," *EURASIP Journal on Advances in Signal Processing*, vol. 1, no. 81, pp. 1–13, 2017.
- [13] J. Qi, J. Du, S.M. Siniscalchi, and C.-H. Lee, "A theory on deep neural network based vector-to-vector regression with an illustration of its expressive power in speech enhancement," *IEEE/ACM Transactions on Audio, Speech, and Language Processing (TASLP)*, vol. 27, no. 12, pp. 1932–1943, 2019.
- [14] D. Yu, F. Seide, G. Li, and L. Deng, "Exploiting sparseness in deep neural networks for large vocabulary speech recognition," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2012, pp. 4409–4412.
- [15] J. Xue, J. Li, and Y. Gong, "Restructuring of deep neural network acoustic models with singular value decomposition," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2013, pp. 2365–2369.
- [16] A. Novikov, D. Podoprikin, A. Osokin, and D. P. Vetrov, "Tensorizing neural networks," in *Advances in Neural Information Processing Systems (NIPS)*, 2015, pp. 442–450.
- [17] I. V. Oseledets, "Tensor-train decomposition," *SIAM Journal on Scientific Computing*, vol. 33, no. 5, pp. 2295–2317, 2011.
- [18] Jun Qi, Xiaoli Ma, Jun Du, Sabato Marco Siniscalchi, and Chin-Hui Lee, "Performance analysis for tensor-train decomposition to deep neural network based vector-to-vector regression," in *54th Annual Conference on Information Sciences and Systems (CISS)*, 2020.
- [19] D. B. P. Douglas and J. M. Baker, "The design for the wall street journal-based csr corpus," in *Proceedings of the workshop on Speech and Natural Language*. Association for Computational Linguistics, 1992, pp. 357–362.
- [20] G. Hu and D. Wang, "A tandem algorithm for pitch estimation and voiced speech segregation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 8, pp. 2067–2079, 2010.
- [21] Eric A Lehmann and Anders M Johansson, "Prediction of energy decay in room impulse responses simulated with an image-source model," *The Journal of the Acoustical Society of America*, vol. 124, no. 1, pp. 269–277, 2008.
- [22] Xiaodi Hou and Liqing Zhang, "Saliency detection: A spectral residual approach," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007, pp. 1–8.
- [23] J. Qi, D. Wang, Y. Jiang, and R. Liu, "Auditory features based on gammatone filters for robust speech recognition," in *2013 IEEE International Symposium on Circuits and Systems (ISCAS2013)*, 2013, pp. 305–308.
- [24] J. Qi, D. Wang, and J. Tejedor Noguerales J. Xu, "Bottleneck features based on gammatone frequency cepstral coefficients," in *Interspeech*. ISCA, 2013, pp. 1751–1755.
- [25] Y. Hirose, K. Yamashita, and S. Hijiya, "Back-propagation algorithm which varies the number of hidden units," *Neural Networks*, vol. 4, no. 1, pp. 61–66, 1991.
- [26] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error log-spectral amplitude estimator," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 33, no. 2, pp. 443–445, 1985.
- [27] J. Qi, D. Wang, and J. Tejedor Noguerales, "Subspace models for bottleneck features," in *Interspeech*. ISCA, 2013, pp. 1746–1750.
- [28] Martín Abadi and et al, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, Software available from tensorflow.org.
- [29] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2001, vol. 2, pp. 749–752.
- [30] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.