

# Modulation spectrum-constrained trajectory error training for mixture density network-based speech synthesis

Sangjun Park, and Minsoo Hahn

Citation: [The Journal of the Acoustical Society of America](#) **144**, EL151 (2018); doi: 10.1121/1.5052206

View online: <https://doi.org/10.1121/1.5052206>

View Table of Contents: <http://asa.scitation.org/toc/jas/144/3>

Published by the [Acoustical Society of America](#)

---

## Articles you may be interested in

[Lower-level acoustics underlie higher-level phonological categories in lexical tone perception](#)

The Journal of the Acoustical Society of America **144**, EL158 (2018); 10.1121/1.5052205

[Computer-based auditory training improves second-language vowel production in spontaneous speech](#)

The Journal of the Acoustical Society of America **144**, EL165 (2018); 10.1121/1.5052201

[Auditory distraction by speech: Comparison of fluctuating and steady speech-like masking sounds](#)

The Journal of the Acoustical Society of America **144**, EL83 (2018); 10.1121/1.5048637

[Perception of relative pitch of sentence-length utterances](#)

The Journal of the Acoustical Society of America **144**, EL89 (2018); 10.1121/1.5048636

[Differences in cue weights for speech perception are correlated for individuals within and across contrasts](#)

The Journal of the Acoustical Society of America **144**, EL172 (2018); 10.1121/1.5052025

[Head waves in ocean acoustic ambient noise: Measurements and modeling](#)

The Journal of the Acoustical Society of America **143**, 1182 (2018); 10.1121/1.5024332

---

# Modulation spectrum-constrained trajectory error training for mixture density network-based speech synthesis

Sangjun Park and Minsoo Hahn<sup>a)</sup>

*School of Electrical Engineering, Korea Advanced Institute of Science and Technology,  
 Daejeon 34141, Republic of Korea  
 psj@kaist.ac.kr, mshahn2@kaist.ac.kr*

**Abstract:** In statistical parametric speech synthesis, a mixture density network is employed to address the limitations of a linear output layer such as pre-computed fixed variances and the unimodal assumption. However, it also has a defect, i.e., it cannot deploy a static-dynamic constraint needed in the training phase for high-quality speech synthesis. To cope with this problem, this paper proposes a training algorithm based on the minimum trajectory error for a mixture density network. And a modulation spectrum-constrained loss function is also proposed to alleviate the over-smoothing effect. The experimental results confirm meaningful improvement both in objective and subjective performance measures.

[DDO]

**Date Received:** July 8, 2018

**Date Accepted:** August 14, 2018

## 1. Introduction

The speech quality of statistical parametric speech synthesis (SPSS) has been improved noticeably through the use of deep neural networks (DNNs),<sup>1</sup> which show better performance in representing the complex, nonlinear, and high-dimensional relationships between linguistic and acoustic features as compared to conventional hidden Markov models.<sup>2</sup> Over the last few years, several end-to-end speech synthesis frameworks generating human-like synthetic speech, including WaveNet, Tacotron, and Deepvoice, have been proposed.<sup>3–5</sup> However, they work so poorly under small corpus conditions and also incur high computation costs. Hence, SPSS approaches are still more useful in real environments.

In SPSS, several studies with various DNN-based architectures have been reported.<sup>6–10</sup> Usually they include a linear output layer and are trained with a mean squared error (MSE) loss function, generating acoustic features with a maximum likelihood parameter generation (MLPG) algorithm.<sup>11</sup> However, two unneglectable error sources still exist. They are the frame-wise independence assumption for the MSE criterion and the unimodal assumption for the linear output layer. Although the temporal information of speech is crucial for high-quality speech synthesis,<sup>12</sup> the MSE criterion fragments the relationship between the static and dynamic features in the training phase. Wu and King proposed the minimum trajectory error (MTE) criterion to moderate this drawback by adding a static-dynamic constraint on the MSE criterion.<sup>13</sup> The MTE training can generate a more natural trajectory, but it is still over-smoothed due to the linear output layer.<sup>14</sup> A mixture density network (MDN) can overcome this problem.<sup>14–17</sup> Multiple Gaussian mixtures in an MDN can represent the multimodality of speech and predict the variances while the linear output layer uses pre-computed fixed variances. Nevertheless, like the MSE criterion, the MDN also generates an unnatural trajectory because it cannot deploy temporal information during the training phase, i.e., the relationships between adjacent frames cannot be used. Another way of alleviating the over-smoothing effect is to utilize analytical features such as the global variance (GV) and the modulation spectrum (MS), known as perceptual cues.<sup>18,19</sup> Training/synthesis algorithms constrained on the GV and MS show improved clarity of synthetic speech.<sup>20–23</sup> However, these methods focus only on making the trajectory variation similar to the natural ones but do not consider the multimodality of speech.

In this paper, a novel MTE criterion-based training algorithm for MDNs is proposed to address both over-smoothing and unnatural trajectory problems. To introduce the MTE criterion, we reformulated the conventional iterative MLPG algorithm for MDNs (Ref. 11) into a closed-form solution utilizing only the most probable mixture (MPM) component. The proposed algorithm covers both the static-dynamic

<sup>a)</sup> Author to whom correspondence should be addressed.

constraint and the multimodality of speech and thus can generate more natural and clear synthetic speech. Furthermore, we introduce an MS constraint into the MTE loss function and mitigate the over-smoothing effect. Our experimental results confirm that the proposed algorithm improves the synthetic speech quality meaningfully both in objective and subjective evaluations.

## 2. Trajectory error training for an MDN constrained on a MS

### 2.1 Conventional training algorithm for a linear output layer

For a given frame-level linguistic feature sequence  $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_t^\top, \mathbf{x}_T^\top]^\top$  and an observed acoustic sequence  $\mathbf{O} = [\mathbf{o}_1^\top, \dots, \mathbf{o}_t^\top, \mathbf{o}_T^\top]^\top$ , the linear output layer in a DNN predicts the acoustic features directly<sup>6-10</sup> as  $\hat{\mathbf{O}} = F_l(\mathbf{X}; \Phi_l)$ , where  $F_l(\cdot)$  is a DNN-based mapping function with a linear output layer and  $\Phi_l$  denotes the parameters of  $F_l(\cdot)$ . Note that the acoustic feature sequence consists of the static and dynamic features as  $\mathbf{o}_t = [\mathbf{c}_t^\top, \Delta \mathbf{c}_t^\top, \Delta^2 \mathbf{c}_t^\top]^\top$ , and  $\mathbf{O}$  can be obtained as  $\mathbf{O} = \mathbf{W}\mathbf{C}$ , where  $\mathbf{W}$  is a weight matrix for calculating dynamic features from static features.<sup>11</sup> To obtain  $\hat{\mathbf{C}}$  from the observed acoustic features  $\hat{\mathbf{O}}$  with the maximum likelihood criterion, an MLPG algorithm is used. It can be written as

$$\hat{\mathbf{C}} = (\mathbf{W}^\top \Sigma^{-1} \mathbf{W})^{-1} \mathbf{W}^\top \Sigma^{-1} \hat{\mathbf{O}}, \quad (1)$$

where  $\Sigma = \text{diag}[\Sigma_1, \dots, \Sigma_t, \dots, \Sigma_T]$  is a covariance matrix sequence.<sup>11</sup> Note that  $F_l(\cdot)$  predicts only the mean of the acoustic features. For this reason,  $\Sigma$  is pre-computed over the training corpus. The MSE loss is designed to minimize the L2 loss of the observed acoustic features  $\hat{\mathbf{O}}$ , whereas the trajectory loss is defined as the L2 loss of the static acoustic features  $\hat{\mathbf{C}}$  in Eq. (1). Now, the MSE and MTE loss functions can be written as follows:

$$L_{\text{MSE}}(\mathbf{O}, \hat{\mathbf{O}}) = \frac{1}{T} (\hat{\mathbf{O}} - \mathbf{O})^\top (\hat{\mathbf{O}} - \mathbf{O}), \quad (2)$$

$$L_{\text{MTE}}(\mathbf{C}, \hat{\mathbf{C}}) = \frac{1}{T} (\hat{\mathbf{C}} - \mathbf{C})^\top (\hat{\mathbf{C}} - \mathbf{C}). \quad (3)$$

The MTE training algorithm has a static-dynamic constraint based on the MLPG algorithm in the training phase and can generate a more natural synthetic speech.<sup>13</sup> However, this method still suffers from over-smoothing due to the unimodal assumption of the linear output layer.

### 2.2 Conventional training algorithm for an MDN

An MDN, which utilizes the multimodal feature distribution using a Gaussian mixture model (GMM), can overcome the above problem of a linear output layer. The MDN maps  $\mathbf{X}$  to a GMM parameter sequence as  $\Lambda = [\lambda_1^\top, \dots, \lambda_t^\top, \dots, \lambda_T^\top]^\top = F_m(\mathbf{X}; \Phi_m)$ ,<sup>14</sup> where  $\lambda_t$  and  $F_m(\cdot)$  denote the GMM parameters at the  $t$ th frame and the DNN-based mapping function with the MDN, respectively.  $F_m(\cdot)$  can be trained by maximizing the GMM likelihood. The maximum likelihood (ML) loss function then becomes

$$L_{\text{ML}}(\mathbf{O}, \Lambda) = -\frac{1}{T} \sum_{t=1}^T \log \sum_{m=1}^M w_{m,t} \cdot N(\mathbf{o}_t; \boldsymbol{\mu}_{m,t}, \boldsymbol{\sigma}_{m,t}), \quad (4)$$

where  $M$  is the number of mixtures and  $w_{m,t}$ ,  $\boldsymbol{\mu}_{m,t}$ , and  $\boldsymbol{\sigma}_{m,t}$  correspond to the mixture weight, the mean, and the variance of the  $m$ th Gaussian component at the  $t$ th frame. From the GMM parameters  $\Lambda$ ,  $\hat{\mathbf{C}}$  can also be obtained using an MLPG algorithm.<sup>11</sup> The MLPG algorithm for the MDN can be written as

$$\mathbf{W}^\top \overline{\Sigma}^{-1} \mathbf{W} \hat{\mathbf{C}} = \mathbf{W}^\top \overline{\Sigma}^{-1} \overline{\mathbf{M}}, \quad (5a)$$

$$\overline{\Sigma}^{-1} = \text{diag}[\overline{\Sigma}_1^{-1}, \dots, \overline{\Sigma}_T^{-1}], \quad (5b)$$

$$\overline{\Sigma}_t^{-1} = \sum_{m=1}^M \gamma_t(m) \overline{\Sigma}_{m,t}^{-1}, \quad (5c)$$

$$\overline{\Sigma}^{-1} \overline{\mathbf{M}} = [\overline{\Sigma}_1^{-1} \boldsymbol{\mu}_1^\top, \dots, \overline{\Sigma}_T^{-1} \boldsymbol{\mu}_T^\top]^\top, \quad (5d)$$

$$\overline{\Sigma}_t^{-1} \boldsymbol{\mu}_t = \sum_{m=1}^M \gamma_t(m) \Sigma_{m,t}^{-1} \boldsymbol{\mu}_{m,t}, \quad (5e)$$

where  $\gamma_t(m)$  is the occupancy probability, which can be solved with an iterative Expectation-Maximization (EM) algorithm. Due to the iterative-form solution, the trajectory loss of MDN cannot be derived, as opposed to Eq. (3).

### 2.3 Proposed trajectory error training algorithm for an MDN

To utilize the aforementioned advantages of both the MTE criterion and an MDN, we propose a new MTE training algorithm for an MDN. To derive the MTE loss function for an MDN, we reformulated Eq. (5a) into a closed-form solution using only the MPM component. This approach can be considered as a suboptimal solution. In this case,  $\gamma_t(m)$  is 1 at every time step with  $M = 1$ . Accordingly, Eqs. (5b) and (5d) can be rewritten as  $\bar{\Sigma}^{-1} = \Sigma_{\text{MPM}}^{-1} = \text{diag}[\Sigma_{m_{\text{MPM},1}^{(1)},1}^{-1}, \dots, \Sigma_{m_{\text{MPM},T}^{(T)},T}^{-1}]$  and  $\bar{\Sigma}^{-1}\bar{M} = \Sigma_{\text{MPM}}^{-1}M_{\text{MPM}} = [\Sigma_{m_{\text{MPM},1}^{(1)},1}^{-1}\mu_{m_{\text{MPM},1}^{(1)},1}, \dots, \Sigma_{m_{\text{MPM},T}^{(T)},T}^{-1}\mu_{m_{\text{MPM},T}^{(T)},T}]^T$ , respectively, where  $m_{\text{MPM}}^{(t)}$  denotes the MPM at the  $t$ th frame. Finally, Eq. (5a) can be reformulated as follows:

$$\begin{aligned} W^T \Sigma_{\text{MPM}}^{-1} W \hat{C} &= W^T \Sigma_{\text{MPM}}^{-1} M_{\text{MPM}}, \\ \hat{C} &= (W^T \Sigma_{\text{MPM}}^{-1} W)^{-1} W^T \Sigma_{\text{MPM}}^{-1} M_{\text{MPM}}. \end{aligned} \quad (6)$$

Using this closed-form solution, the MTE loss function for the MDN can be defined as follows:

$$\begin{aligned} L_{\text{MTE2}}(C, \hat{C}) &= \frac{1}{T} \left( (W^T \Sigma_{\text{MPM}}^{-1} W)^{-1} W^T \Sigma_{\text{MPM}}^{-1} M_{\text{MPM}} - C \right)^T \\ &\quad \times \left( (W^T \Sigma_{\text{MPM}}^{-1} W)^{-1} W^T \Sigma_{\text{MPM}}^{-1} M_{\text{MPM}} - C \right). \end{aligned} \quad (7)$$

Note that  $\Sigma$  is also trainable, in contrast to Eq. (3). By the way, two ways to determine the MPM can be expressed as follows:

$$m_{\text{MPM}}^{(t)} = \arg \max_m w_{m,t}, \quad (8)$$

$$m_{\text{MPM}}^{(t)} = \arg \max_m N(\mathbf{o}_t; \mu_{m,t}, \sigma_{m,t}). \quad (9)$$

In the synthesis phase, the MPM is determined by Eq. (8) because  $\mathbf{o}_t$  in Eq. (9) is unobservable. On the other hand, in order to train the mixtures corresponding to the training data, the MPM in Eq. (9) should be used in the training phase and it then can be considered as a teacher-forcing method.<sup>24</sup> If the MPM in Eq. (8) is adopted in the training phase, only one mixture is chosen as the MPM for the given linguistic feature, and the output of the DNN can be smoothed similarly to a linear output layer. In addition, when the MDN is trained by only the MTE loss function in Eq. (7), the parameters of the mixtures except for those of the MPM would be updated erroneously because they are affected by the shared layers. In order to train the MDN stably, we propose a new loss function to optimize Eqs. (4) and (7) jointly, as follows:

$$L_{\text{MTE MDN}}(C, \hat{C}) = L_{\text{ML}}(WC, \Lambda) + L_{\text{MTE2}}(C, \hat{C}). \quad (10)$$

### 2.4 Proposed MS-constrained training algorithm

Although the training algorithm in Sec. 2.3 can generate a more natural trajectory due to the use of a static-dynamic constraint, the generated features can still be smoothed by outliers because they have not a multimodal but an irregular distribution. To address this problem, we propose an MS-constrained MTE loss function. It improves the synthetic speech quality by making the trajectory variation similar to natural ones. We defined the MS as a log power spectrum of the acoustic feature sequence<sup>21</sup>  $\mathbf{S}(C) = [\mathbf{S}_1, \dots, \mathbf{S}_D]$ ,  $\mathbf{S}_d = [\mathbf{S}_d^0, \dots, \mathbf{S}_d^N]$ , where  $D$  and  $N$  correspondingly denote the acoustic feature dimension and half of the fast Fourier transform (FFT) size. The segment-level MS, not the utterance-level MS, is used due to the normalization effect; it is computed regardless of the utterance length without zero-padding.<sup>21</sup> The MS loss function can be written as

$$L_{\text{MS}}(C, \hat{C}) = \frac{1}{K} (\mathbf{S}(\hat{C}) - \mathbf{S}(C))^T (\mathbf{S}(\hat{C}) - \mathbf{S}(C)), \quad (11)$$

where  $K$ ,  $\mathbf{S}(C)$ , and  $\mathbf{S}(\hat{C})$  are the number of MS segments and the MS of natural and generated acoustic features, respectively. From Eqs. (10) and (11), the MS-constrained MTE loss function for the MDN is defined as

Table 1. Acoustic model configurations for each experimental system.

Notation	Output layer	Training phase	Synthesis phase
MTE	Linear output layer	Update $\Phi_l$ with Eq. (3)	Generate $\hat{C}$ with Eq. (1)
ML	MDN output layer	Update $\Phi_m$ with Eq. (4)	Generate $\hat{C}$ with Eq. (5a)
Proposed-1	MDN output layer	Update $\Phi_m$ with Eqs. (9) and (10)	Generate $\hat{C}$ with Eqs. (6) and (8)
Proposed-2	MDN output layer	Update $\Phi_m$ with Eqs. (9) and (13)	Generate $\hat{C}$ with Eqs. (6) and (8)

$$L_{\text{MTE MDN MS}}(\mathbf{C}, \hat{\mathbf{C}}) = L_{\text{MTE MDN}}(\mathbf{C}, \hat{\mathbf{C}}) + L_{\text{MS}}(\mathbf{C}, \hat{\mathbf{C}}). \quad (12)$$

In contrast to Eq. (10), each term in Eq. (12) has a different perspective on improving the speech quality;  $L_{\text{MTE MDN}}(\mathbf{C}, \hat{\mathbf{C}})$  and  $L_{\text{MS}}(\mathbf{C}, \hat{\mathbf{C}})$  involve the generation loss and variation of speech, respectively. If the momentum is biased for  $L_{\text{MS}}(\mathbf{C}, \hat{\mathbf{C}})$ , the trained model focuses only on the high-frequency component of the features and not on the generation loss. This problem can be solved by adding an emphasis coefficient  $\alpha$ , which is used in an MS-based postfilter.<sup>21</sup> The final loss function can be written as follows:

$$L_{\text{MTE MDN MS}}(\mathbf{C}, \hat{\mathbf{C}}) = (1 - \alpha)L_{\text{MTE MDN}}(\mathbf{C}, \hat{\mathbf{C}}) + \alpha L_{\text{MS}}(\mathbf{C}, \hat{\mathbf{C}}). \quad (13)$$

Hence, a larger  $\alpha$  is for a more natural variation in speech, while a smaller  $\alpha$ , for less generation loss in training.

### 3. Experiments and results

#### 3.1 Experiment configuration

The Blizzard Challenge 2013 Nancy corpus was used as the training data. It consists of 12 095 English utterances ( $\sim 18$  h) sampled at 16 kHz.<sup>25</sup> 90% and 10% of the corpus except for 100 evaluation utterances were used as the training and validation data, respectively. The WORLD vocoder was utilized for the analysis/synthesis of the acoustic features.<sup>26</sup> The 25th order Mel-cepstrum coefficients (MCs), coded aperiodicity (AP), continuous  $F_0$  (Ref. 27), and voiced/unvoiced (V/UV) measures were extracted for each frame with a 5 ms shift. In addition, 372 linguistic features, e.g., the dot level, stress, inflection, playability, and the quin-Lessemes identity, were extracted based on Lessemes.<sup>28</sup>

To investigate the performance of the proposed algorithm, experimental systems were constructed by adopting a bi-directional long short-term memory (bLSTM) recurrent neural network for a duration model (DM) and an acoustic model (AM). Three bLSTM layers with 64 and 256 memory cells were used for the DM and the AM, respectively. The linear output layer trained with the MSE criterion by Eq. (2) was used in the DM. The detailed AM configurations are given in Table 1. Empirically determined different numbers of mixtures for each acoustic feature were used; four for MC, two for AP, two for  $F_0$ , and one for V/UV. The MS is computed from the acoustic feature sequence but without V/UV. For segment-level MS computations, a Bartlett window with a frame-length of 25 and with a frame shift size of 12 was used, and the FFT size was 64 frames. The emphasis coefficient  $\alpha$  of Proposed-2 was determined empirically through listening tests. We set  $\alpha$  to 0.2, because the MTE loss increases rapidly when  $\alpha \geq 0.3$  and we also found some unstable fluctuations for some samples when  $\alpha = 0.3$ . All systems were trained using a mini-batch stochastic gradient descent-based backpropagation algorithm with the Adam optimizer.<sup>29</sup> The early stopping method was adopted to determine the number of training epochs.<sup>30</sup> All systems were implemented in PyTorch.<sup>31</sup>

For a subjective evaluation, preference tests of four cases were conducted; 20 listeners assessed 20 synthetic speech pairs for each case shown in Table 1; the Mel-cepstral distortion (MCD) and root-mean-square error (RMSE) of  $F_0$  were computed

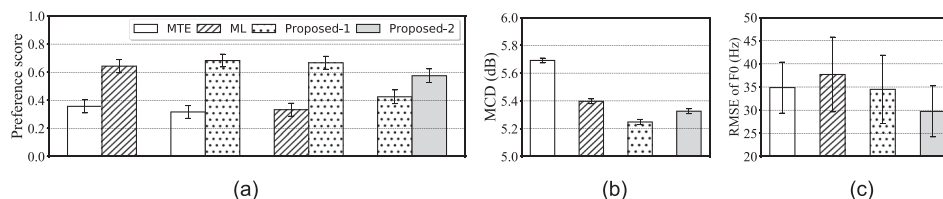


Fig. 1. Subjective and objective evaluation results with a 95% confidence interval: (a) preference score for four cases: MTE vs ML, MTE vs Proposed-1, ML vs Proposed-1 and Proposed-1 vs Proposed-2 in order; (b) MCD; and (c) RMSE of  $F_0$ .



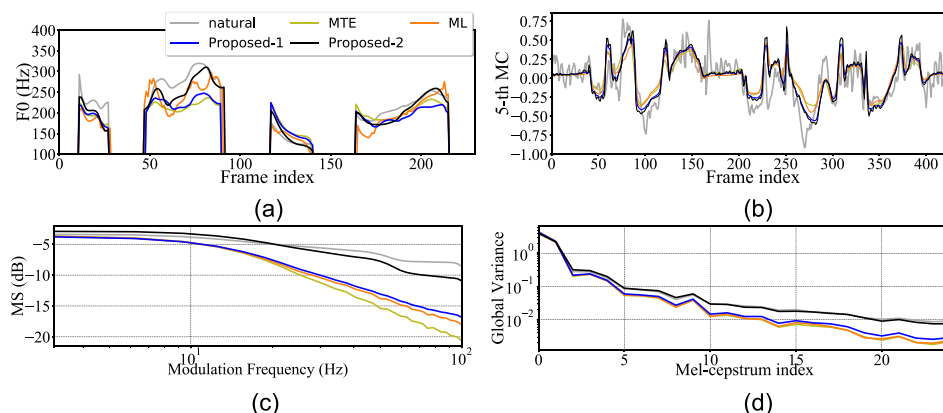


Fig. 2. (Color online) Comparative plot of the generated features with the natural features in  $F_0$ , MC, MS, and GV-domains: (a)  $F_0$  contour example; (b) fifth MC contour example; (c) averaged MS of the 20th MC sequences; and (d) averaged GV of MC sequences.

for an objective evaluation. The speech samples were generated from the DM and the AM for the subjective test, whereas they were generated from only the AMs for the objective test. All experimental results were drawn for 100 evaluation utterances.

### 3.2 Results and discussion

Figure 1 summarizes our subjective and objective test results. Figure 2 shows the generated feature plot with its natural counterpart. ML and the proposed methods are preferred to MTE, as shown in Fig. 1(a). Although MTE predicted  $F_0$  fairly well by deploying temporal information, Fig. 1(b) shows that it is limited when used to model the MC multimodality, as described in earlier work.<sup>15</sup> Figure 2(b) also confirms that the generated MC sequence from MTE is smoother than the others.

Proposed-1 outperforms ML on the objective and subjective measures. Specifically,  $F_0$  generated from ML is degraded by abrupt fluctuations in time, as illustrated in Figs. 1(c) and 2(a). This phenomenon also arose in earlier work,<sup>17</sup> having been caused by the absence of dynamic statistics in the training procedure. By introducing the MTE criterion, Proposed-1 moderated this drawback of the ML criterion. However, the features generated by Proposed-1 are still smoothed by the outliers.

By introducing an MS constraint, Proposed-2 improves this smoothing problem with the better MS and GV, which are more similar to the natural ones, as shown in Figs. 2(c) and 2(d). Note that the lower MS and GV mean a more smoothed and degraded trajectory. A significant improvement in the  $F_0$  estimation was achieved with minor degradation in the MC estimation, as depicted in Figs. 1(b) and 1(c). We can therefore conclude that certain amounts of natural fluctuations of  $F_0$  and the MC in the time domain improve the naturalness and clarity of synthetic speech despite the higher generation loss. Our experimental result also supports this idea. Whereas the MC values generated by other algorithms were narrowly distributed, those by Proposed-2 were rather widely distributed like that of the natural samples.

## 4. Conclusion

In this paper, a novel training algorithm for an MDN based on the MTE criterion was proposed in an attempt to overcome the quality degradation of synthetic speech caused by the frame-wise independence and unimodal assumptions of conventional DNN-based SPSS algorithms.

Our proposed algorithm improved the naturalness and the clarity of synthetic speech with more than a 60% preference score in a subjective evaluation when compared to the conventional algorithms and more precise  $F_0$  and unsmoothed MC in an objective evaluation. We also investigated an MS constraint which reduces over-smoothing caused by outliers and got a reasonably successful improvement. Furthermore, the proposed parameter generation algorithm in a closed-form can save the computation cost of the iterative MLPG algorithm. Considering all these results, we can say that our proposed methods achieved a meaningful improvement in the synthetic speech quality.

### Acknowledgments

This research was supported by the Basic Science Research Program through the National Research Foundation of Korea funded by the Ministry of Science, ICT, and Future Planning (Grant No. 2017R1A2B4011357).

## References and links

- <sup>1</sup>Z. H. Ling, S.-Y. Kang, H. Zen, A. Senior, M. Schuster, X.-J. Qian, H. M. Meng, and L. Deng, “Deep learning for acoustic modeling in parametric speech generation: A systematic review of existing techniques and future trends,” *IEEE Signal Process. Mag.* **32**(3), 35–52 (2015).
- <sup>2</sup>K. Tokuda, Y. Nankaku, T. Toda, H. Zen, J. Yamagishi, and K. Oura, “Speech synthesis based on hidden Markov models,” *Proc. IEEE* **101**(5), 1234–1252 (2013).
- <sup>3</sup>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,” CoRR abs/1609.03499 (2016), [arxiv.org/abs/1609.03499](https://arxiv.org/abs/1609.03499).
- <sup>4</sup>J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. J. Skerry-Ryan, R. A. Saurous, Y. Agiomyrgiannakis, and Y. Wu, “Natural TTS synthesis by conditioning Wavenet on mel spectrogram predictions,” in *Proceedings of the 2018 IEEE ICASSP*, Calgary, Canada (2018), pp. 4779–4783.
- <sup>5</sup>W. Ping, K. Peng, A. Gibiansky, S. O. Arik, A. Kannan, S. Narang, J. Raiman, and J. Miller, “Deep voice 3: Scaling text-to-speech with convolutional sequence learning,” in *Proceedings of the International Conference on Learning Representations (ICLR 2018)*, Vancouver, Canada (2018).
- <sup>6</sup>H. Zen, A. Senior, and M. Schuster, “Statistical parametric speech synthesis using deep neural networks,” in *Proceedings of the 2013 IEEE ICASSP*, Vancouver, Canada (2013), pp. 7962–7966.
- <sup>7</sup>Z. H. Ling, L. Deng, and D. Yu, “Modeling spectral envelopes using restricted Boltzmann machines and deep belief networks for statistical parametric speech synthesis,” *IEEE Trans. Audio, Speech, Lang. Process.* **21**(10), 2129–2139 (2013).
- <sup>8</sup>H. Zen and H. Sak, “Unidirectional long short-term memory recurrent neural network with recurrent output layer for low-latency speech synthesis,” in *Proceedings of the 2015 IEEE ICASSP*, Brisbane, Australia (2015), pp. 4470–4474.
- <sup>9</sup>Z. Wu, C. Valentini-Botinhao, O. Watts, and S. King, “Deep neural networks employing multi-task learning and stacked bottleneck features for speech synthesis,” in *Proceedings of the 2015 IEEE ICASSP*, Brisbane, Australia (2015), pp. 4460–4464.
- <sup>10</sup>H. Zen, Y. Agiomyrgiannakis, N. Egberts, F. Henderson, and P. Szczepaniak, “Fast, compact, and high quality LSTM-RNN based statistical parametric speech synthesizer for mobile devices,” in *Proceedings of INTERSPEECH 2016*, San Francisco, CA (2016), pp. 2273–2277.
- <sup>11</sup>K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” in *Proceedings of the 2000 IEEE ICASSP* (2000), Vol. 3, pp. 1315–1318.
- <sup>12</sup>L. Xu and Y. Zheng, “Spectral and temporal cues for phoneme recognition in noise,” *J. Acoust. Soc. Am.* **122**(3), 1758–1764 (2007).
- <sup>13</sup>Z. Wu and S. King, “Improving trajectory modelling for DNN-based speech synthesis by using stacked bottleneck features and minimum generation error training,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.* **24**(7), 1255–1265 (2016).
- <sup>14</sup>C. M. Bishop, “Mixture density networks,” Technical Report, Birmingham, AL, [publications.aston.ac.uk/373/](http://publications.aston.ac.uk/373/).
- <sup>15</sup>H. Zen and A. Senior, “Deep mixture density networks for acoustic modeling in statistical parametric speech synthesis,” in *Proceedings of the 2014 IEEE ICASSP*, Florence, Italy (2014), pp. 3844–3848.
- <sup>16</sup>K. Richmond, “Trajectory mixture density networks with multiple mixtures for acoustic-articulatory inversion,” in *Advances in Nonlinear Speech Processing* (Springer, Berlin, Heidelberg, 2007), pp. 263–272.
- <sup>17</sup>X. Wang, S. Takaki, and J. Yamagishi, “An autoregressive recurrent mixture density network for parametric speech synthesis,” in *Proceedings of the 2017 IEEE ICASSP*, New Orleans, LA (2017), pp. 4895–4899.
- <sup>18</sup>V. Tyagi, I. McCowan, H. Misra, and H. Bourlard, “Mel-cepstrum modulation spectrum (mcms) features for robust ASR,” in *2003 IEEE Workshop on ASRU* (2003), pp. 399–404.
- <sup>19</sup>R. Drullman, J. M. Festen, and R. Plomp, “Effect of reducing slow temporal modulations on speech reception,” *J. Acoust. Soc. Am.* **95**(5), 2670–2680 (1994).
- <sup>20</sup>T. Toda and K. Tokuda, “A speech parameter generation algorithm considering global variance for HMM-based speech synthesis,” *IEICE Trans. Inf. Syst.* **E90-D**(5), 816–824 (2007).
- <sup>21</sup>S. Takamichi, T. Toda, A. W. Black, G. Neubig, S. Sakti, and S. Nakamura, “Postfilters to modify the modulation spectrum for statistical parametric speech synthesis,” *IEEE Trans. Audio, Speech, Lang. Process.* **24**(4), 755–767 (2016).
- <sup>22</sup>S. Takamichi, T. Toda, A. W. Black, and S. Nakamura, “Modulation spectrum-constrained trajectory training algorithm for GMM-based voice conversion,” in *Proceedings of the 2015 IEEE ICASSP*, Brisbane, Australia (2015), pp. 4859–4863.
- <sup>23</sup>K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda, “Trajectory training considering global variance for speech synthesis based on neural networks,” in *Proceedings of the 2016 IEEE ICASSP*, Shanghai, China (2016), pp. 5600–5604.
- <sup>24</sup>R. J. Williams and D. Zipser, “A learning algorithm for continually running fully recurrent neural networks,” *Neural Comput.* **1**(2), 270–280 (1989).
- <sup>25</sup>S. King, K. T. A. W. Black, and K. Prahallad, “The blizzard challenge 2013,” in *Proceedings of the Blizzard Challenge 2013*, Barcelona, Spain (2013), pp. 1–10.
- <sup>26</sup>M. Morise, F. Yokomori, and K. Ozawa, “WORLD: A vocoder-based high-quality speech synthesis system for real-time applications,” *IEICE Trans. Inf. Syst.* **E99-D**(7), 1877–1884 (2016).
- <sup>27</sup>K. Yu and S. Young, “Continuous F0 modeling for HMM-based statistical parametric speech synthesis,” *IEEE Trans. Audio, Speech, Lang. Process.* **19**(5), 1071–1079 (2011).

- <sup>28</sup>M. Munro, S. Turner, A. Munro, and K. Campbell, "Use of Lesseemes in text-to-speech synthesis," in *Collective Writings on the Lessac Voice and Body Work: A Festschrift* (Llumina Press, Coral Springs, FL, 2009), pp. 362–374.
- <sup>29</sup>D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," CoRR abs/1412.6980 (2014), [arxiv.org/abs/1412.6980](https://arxiv.org/abs/1412.6980).
- <sup>30</sup>R. Caruana, S. Lawrence, and L. Giles, "Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping," in *Advances in Neural Information Processing Systems* (MIT Press, Cambridge, MA, 2001), pp. 402–408.
- <sup>31</sup>A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in *Proceedings of the NIPS Autodiff Workshop* (2017).