

ENHANCEMENT OF CODED SPEECH USING A MASK-BASED POST-FILTER

Srikanth Korse¹, Kishan Gupta^{1,2} and Guillaume Fuchs^{1,2}

¹Fraunhofer IIS, Erlangen, Germany, srikanth.korse@iis.fraunhofer.de

²International Audio Laboratories, Friedrich-Alexander University (FAU), Erlangen, Germany *

ABSTRACT

The quality of speech codecs deteriorates at low bitrates due to high quantization noise. A post-filter is generally employed to enhance the quality of the coded speech. In this paper, a data-driven post-filter relying on masking in the time-frequency domain is proposed. A fully connected neural network (FCNN), a convolutional encoder-decoder (CED) network and a long short-term memory (LSTM) network are implemented to estimate a real-valued mask per time-frequency bin. The proposed models were tested on the five lowest operating modes (6.65 kbps-15.85 kbps) of the Adaptive Multi-Rate Wideband codec (AMR-WB). Both objective and subjective evaluations confirm the enhancement of the coded speech and also show the superiority of the mask-based neural network system over a conventional heuristic post-filter used in the standard like ITU-T G.718.

Index Terms: Deep Neural Network, Speech Coding, Speech Coding Enhancement, Post-filtering

1. INTRODUCTION

State-of-the-art communication speech codecs such as 3GPP Adaptive Multi-Rate Wide-Band (AMR-WB) [1] and 3GPP Enhanced Voice Services (EVS) [2] use Code-Excited Linear Prediction (CELP) as a core coder for coding speech. CELP consists of 3 essential parts: a short-term prediction using linear predictive coding (LPC), a long-term prediction (LTP) exploiting the fundamental frequency and the innovative codebook for modeling the residual of the predictions. At moderate to high bitrates, sufficient bits are assigned to the LPC coefficients, LTP parameters and the innovative codebook, thereby yielding sufficiently good to transparent quality. However, at low-bitrates, the majority of available bits are allocated to LPC and LTP parameters, and very few bits are left for the innovative codebook. This significantly affects the quality of coded speech at low bit rates.

To enhance the perceptual quality at low bitrates, modern speech codecs employ post-filters [3, 4, 5, 6]. Post-processing the coded speech is conceptually similar to most speech enhancement techniques aiming to attenuate or amplify frequency components with a bad or good signal-to-noise ratio, respectively. However, instead of modelling the noise or the speech, they rely mainly on prior knowledge and parameters of the coding scheme. Typical examples of post-processing are post-filters that emphasize the formants and the pitch structures of the coded speech by reusing information from LPC or LTP, respectively [3, 4]. It is also noteworthy that in [7] a statistical model of the quantization noise of a very simplistic coding scheme is derived for designing an optimal post-filter in the log-magnitude domain.

In recent years, data-driven approaches for speech enhancement and dereverberation have been shown to outperform classical statistical signal processing approaches [8, 9, 10, 11]. They usually rely on data and make no assumptions about the signal or noise statistics. Among them, mask-based approaches [10, 11], estimating either a real-valued ideal ratio mask (IRM), an ideal binary mask (IBM) or a complex-valued ideal ratio mask (cIRM) in the frequency domain, are especially efficient for denoising or inverse filtering tasks. These mask-based approaches have not been employed until now to enhance the quality of the coded speech and the authors hence investigate the benefit of using these approaches in the context of speech coding.

Alternatively, the DNN can also learn the spectral mapping function between the spectral coefficients of noisy or reverberated speech and those of the clean speech [8, 9]. This concept was recently adopted in [12, 13] for enhancing the coded speech in the cepstral domain. In order to reduce the complexity, the cepstrum is truncated and then fed as input to the neural network.

Recently, autoregressive models such as WaveNet [14] and LPC-Net [15] have been employed to enhance the coded speech [16]. The high delay and complexity of these models prevent using them for real-time communication [17]. Therefore, we do not consider these models for our evaluation.

1.1. Key Contribution of this Paper

- The paper shows that a mask based post-filter in the spectral domain performs better than cepstral-domain post-filter (Cepstrum-CNN) as proposed in [12, 13].
- FCNN, CED and LSTM structures were implemented to estimate a real-valued mask per time-frequency bin. This mask was then used to enhance the perceptual quality of the coded speech.
- Among the three models proposed, CED performs better than FCNN and LSTM. Hence, CED was used for the final evaluation.
- The proposed model was trained on the coded speech obtained from the AMR-WB codec at bitrate 6.65 kbps. The trained model was then tested on all bitrates ranging from 6.65 to 15.85 kbps.
- Robustness of the trained model was validated by testing on a completely different database.
- The proposed model is also compared with the heuristic post-filter adopted in G.718 [4].

*International Audio Laboratories is a joint institution between Fraunhofer IIS and Friedrich-Alexander University (FAU)

Mask Thresholds	6.65kbps	8.85kbps	12.65kbps
[0, 1]	38.94%	41.00%	44.09%
(1, 2]	31.19%	33.44%	36.20%
(2, 5]	21.40%	18.69%	14.66%
(5, ∞]	8.46%	6.87%	5.05%

Table 1. Percentage of real-valued mask in different threshold regions measured at lowest three birates of AMR-WB.

2. PROBLEM FORMULATION

2.1. Oracle Experiments

From a simplistic mathematical point of view, one can describe the coded speech $\tilde{x}(n)$ as:

$$\tilde{x}(n) = x(n) + \delta(n) \quad (1)$$

where $x(n)$ is the input speech to the codec and $\delta(n)$ is the quantization noise. The quantization noise $\delta(n)$ is correlated with the input speech since CELP uses a perceptual model during the quantization process. The correlation of quantization noise with the input speech makes our post-filtering problem unique to speech enhancement problem which usually assumes the noise to be uncorrelated.

In order to reduce the quantization noise, we estimate a real-valued mask per time-frequency bin and multiply this mask with the corresponding spectral magnitude of the coded speech as shown:

$$|\hat{X}(k, n)| = M(k, n) \cdot |\tilde{X}(k, n)|, \quad (2)$$

where $M(k, n)$ is the real-valued mask, $|\tilde{X}(k, n)|$ is the magnitude spectrum of the coded speech, $|\hat{X}(k, n)|$ is the magnitude spectrum of the enhanced speech, k is the frequency index and n is the time index.

If the mask is ideal, we can get back the clean speech magnitude spectrum from the coded speech magnitude spectrum. The definition of the ideal ratio mask (IRM) is given by:

$$\text{IRM}(k, n) = \frac{|X(k, n)|}{|\tilde{X}(k, n)| + \gamma}, \quad (3)$$

where $|X(k, n)|$ is the magnitude spectrum of the clean speech and γ is a very small constant factor to prevent division by zero. Since the magnitude values lie in the range 0 to ∞, the values of IRM is positive and unbounded.

Table 1 shows the distribution of real-valued mask per time-frequency bin in different threshold regions at lowest three birates of AMR-WB. These mask values were computed using (3). From the Table 1, it can be concluded that, since most the mask values lie in the region between 0 and 5, we cannot afford to bound the mask value to 1 as usually done in the conventional speech enhancement techniques.

To find out the ideal bounding value and also the performance of the mask-based technique over cepstral methods, oracle experiments were performed at 6.65 kbps and 8.85 kbps. The oracle mask was computed using (3) and then was bounded to values 1, 2, 4 and 10. The oracle cepstrum was computed as follows: First, cepstrum of length 512 was obtained from both coded and clean speech. Second, the first 64 cepstral values obtained from the coded speech were replaced by the clean speech cepstral values.

Fig. 1 compares the average Perceptual Objective Listening Quality Assessment (POLQA) [18] scores between coded, oracle cepstrum and oracle mask with different bounding values at birates

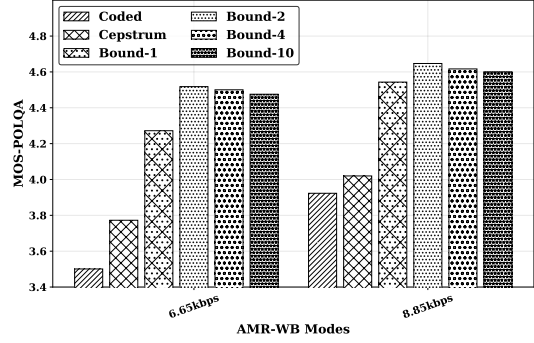


Fig. 1. Average POLQA scores evaluating the oracle experiment

6.65 kbps and 8.85 kbps. It can be clearly observed that, in the oracle case, spectral mask-based methods outperform the cepstrum method. With regards to the upper bound, it can be observed that the bound values greater than 1 perform significantly better than bound 1. Among the bound values 2, 4 and 10, there is no significant difference. This motivated us to bound our mask to 2 in further experiments.

The above observations in the oracle experiments motivate the adoption of mask-based approach in the spectral domain for enhancing the perceptual quality of the coded speech.

2.2. Modified Signal Approximation

We train the neural network using modified signal approximation (mod-SA). The major differences between mod-SA and traditional signal approximation (SA) defined in [11] are as follows:

- Instead of computing the IRM and bounding it to 2, the modified mask $\tilde{M}(k, n)$ is computed as:

$$\tilde{M}(k, n) = \begin{cases} \text{IRM}(k, n) & \text{if } \text{IRM}(k, n) \leq \alpha \\ \rho & \text{if } \text{IRM}(k, n) > \alpha \end{cases} \quad (4)$$

where $0 \leq \rho \leq \alpha$ and $\text{IRM}(k, n)$ is computed using (3). For our experiments, we choose α as 2 and ρ as 1. In other words, when the IRM is greater than bound 2, the coded spectral magnitude is kept unchanged.

- Since the mask was modified, the target is also modified and the modified target $|\bar{X}(k, n)|$ is given by:

$$|\bar{X}(k, n)| = \tilde{M}(k, n) \cdot |\tilde{X}(k, n)|, \quad (5)$$

- The mean square error (MSE) is computed between the modified target and enhanced speech in the log-magnitude domain instead of the magnitude domain.

The above modifications are essential to get a generalized model that works on the higher bitrates despite being trained only on the lowest one. If the objective is to have different models for different bitrates, the authors would advise to use then the unmodified target $|X(k, n)|$ for training.

Layer name	Input	Hyperparameter	Output
Reshape	6 × 205	-	1 × 6 × 205
Conv2d.1	1 × 6 × 205	2 × 3, (1,2), 16	16 × 5 × 102
Conv2d.2	16 × 5 × 102	2 × 3, (1,2), 32	32 × 4 × 50
Conv2d.3	32 × 4 × 50	2 × 3, (1,2), 64	64 × 3 × 24
Conv2d.4	64 × 3 × 24	2 × 3, (1,2), 128	128 × 2 × 11
Deconv2d.1	128 × 2 × 11	2 × 3, (1,2), 64	64 × 3 × 23
Deconv2d.2	128 × 3 × 24	2 × 3, (1,2), 32	32 × 4 × 49
Deconv2d.3	64 × 3 × 50	2 × 3, (1,2), 16	16 × 5 × 101
Deconv2d.4	32 × 5 × 102	2 × 3, (1,2), 1	1 × 6 × 205
Conv2d.5	1 × 6 × 205	6 × 1, (1,1), 1	1 × 1 × 205
Flatten	1 × 6 × 205	-	1 × 205

Table 2. Architecture of our proposed CED. The input and output size is given as $\text{featureMaps} \times \text{timesteps} \times \text{frequencyBins}$. The hyperparameter is indicated as kernelSize , strides , outchannels .

3. EXPERIMENTAL SETUP

Our proposed post-filter computes the short-time Fourier transform (STFT) of frames of 32 ms with 50% overlap (16 ms) at 16 kHz sampling rate resulting in 257 frequency bins. The square root of the Hann window is used as analysis and synthesis window. Only bandwidth up to 6.4 kHz (205 frequency bins) was processed with DNN and the frequency region between 6.4 to 7 kHz was left unprocessed. Since speech has temporal dependency, past frames were used as context frames i.e. input to the DNN were past frames plus the current frame. The input to the DNN was normalized log-magnitude since magnitude values have a higher dynamic range compared to log-magnitude. The output of the DNN is a real-valued mask that lies in the range 0 to 2. The real-valued mask is then multiplied with the coded magnitude to obtain the enhanced magnitude.

The three proposed networks are explained below in detail:

- **FCNN:** Three past frames and the current frame are concatenated and provided as input to the FCNN. The size of the input was 820. It has two hidden layers with 1024 units. Each hidden layer consists of Rectified Linear Units (ReLU) as activation functions along with batch normalization and a dropout of 0.2. The output layer consists of 205 units.
- **LSTM:** The LSTM network consists of two LSTM layers with 400 and 205 units, respectively, with 10 time steps (9 past frames and current frame). A Dropout of 0.1 and recurrent dropout of 0.2 was used. Only the last time step of the second LSTM unit was given as input to the output layer.
- **CED:** An encoder-decoder architecture-based CNN is implemented as shown in Table 2. The input to the CED is 6 time steps (5 past frames and current frame). Each layer of CNN uses batch normalization and ELU (Exponential Linear Unit) activation function. Skip connections are used between encoder and decoder. Required zero-padding is done in the time frame to match the frequencyBins dimensions for skip connections.

For all the models, the output layer consists of sigmoid units scaled with factor 2. All the models were trained with the ADAM optimizer [19] with a learning rate of 0.001 and a batch size of 32. Instead of using a fixed number of epochs, training was done till convergence using early stopping. The implementation of the reference Cepstrum-CNN is the same as proposed in [12]. It is to be noted that the Cepstrum-CNN do not use any context frames.

Table 3 shows the number of parameters and frame sizes for each of the compared networks. Our proposed CED model has the

Network Architecture	Number of Parameters	Frame Size
FCNN	2108621	32ms
LSTM	1468120	32ms
CED	147292	32ms
Cepstrum-CNN	419805	20ms

Table 3. Comparison of the number of parameters in different network architectures.

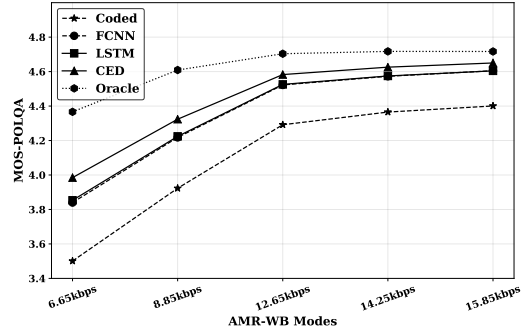


Fig. 2. POLQA scores evaluating the performance of the FCNN, LSTM and CED architectures using the NTT test set.

smallest number of parameters in comparison to all other models. On other hand, our model operates on 32 ms frames compared to Cepstrum-CNN which operates on 20 ms. This results in 16 ms delay for our model instead of 10 ms delay for Cepstrum-CNN, which is still acceptable for real time communication.

The input speech signals were first filtered using the P.341 filter (cutoff frequency of 7kHz) [20] and then the active speech level was adjusted to -26 dBov [21] before coding with AMR-WB. The coded speech signals were also filtered using P.341 filter with same 7 kHz cutoff frequency before testing. For objective assessment, we used POLQA, while for subjective assessment, we followed the methodology Multiple Stimuli with Hidden Reference and Anchor (MUSHRA) [22].

The time-domain enhanced speech was obtained using the inverse short-time Fourier transform (iSTFT) and a synthesis window before the overlap and add step. iSTFT made use of the phase of the coded speech without any processing.

4. EXPERIMENTS AND RESULTS

For training, we used the NTT-AT [23] database. The files were downsampled to 16 kHz and a passive mono downmix was obtained from the stereo files. Out of 3690 files, 3612 files were used for training, 198 files were used for validation and 150 files were used for testing. The database was split in such a way that the distribution in terms of male and female speakers and languages was balanced in the training, validation and test set. All the files were preprocessed as explained in Section 3.

Fig. 2 compares the POLQA scores of the three proposed architectures (FCNN, CED and LSTM). Among the proposed architectures, LSTM and FCNN have the same performance while CED consistently performs better than the two others across all bitrates. Hence, for further evaluation, we only consider the CED architecture. The oracle mask used for comparison was obtained as explained in Section 2.2.

Fig. 3 compares the POLQA scores of the proposed CED model with the post-filter used in G.718 and the Cepstrum-CNN model.

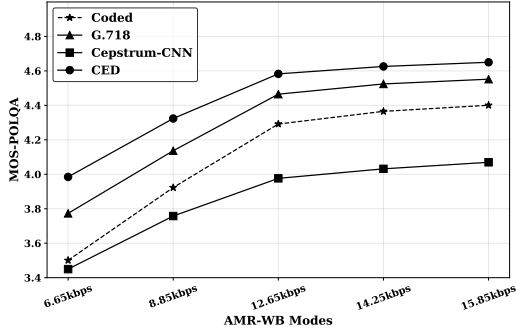


Fig. 3. POLQA scores evaluating the performance of the Cepstrum-CNN, CED and G.718 using the NTT test set.

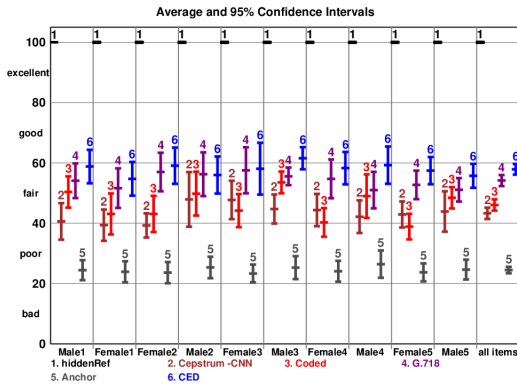


Fig. 4. Average MUSHRA scores of 11 listeners at 6.65 kbps.

Our proposed CED model improves the perceptual quality of the coded speech and is consistently better than G.718 at all bitrates. The Cepstrum-CNN fails to improve the perceptual quality of the coded speech at all bitrates.

The POLQA scores are consistent with the MUSHRA scores as shown in Fig. 4 and Fig. 5 which compare the proposed CED model with G.718 and the Cepstrum-CNN at bitrates 6.65 and 12.65 kbps, respectively. Both listening tests involve 11 expert listeners. The listening tests were conducted in the listening test rooms which were isolated from the outside noise. STAX headphones were used for the listening tests.

At 6.65 kbps, our proposed CED model successfully enhances the coded speech and is also better than the G.718. Compared to the coded speech, the CED model gains around 12 MUSHRA points and 0.5 POLQA mean opinion score (MOS). At 12.65 kbps, the benefit of our proposed CED model is less and is close to G.718 post-filter. This is primarily because the quality of the coded speech lies already in the good range.

The Cepstrum-CNN model fails to enhance the quality of coded speech at both 6.65 kbps and 12.65 kbps. This is due to the following reasons:

- Since only the first cepstral coefficients are enhanced, it affects only the spectral envelope and not the spectral fine structures.
- The Cepstrum-CNN fails to suppress the noise between harmonics compared to our proposed CED model or the G.718 post-filter.
- At high frequencies, the Cepstrum-CNN is closer to the target

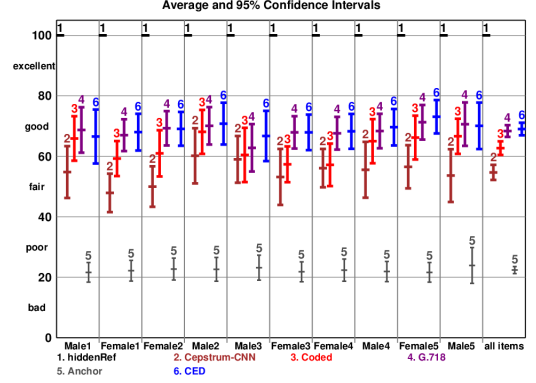


Fig. 5. Average MUSHRA scores of 11 listeners at 12.65 kbps.

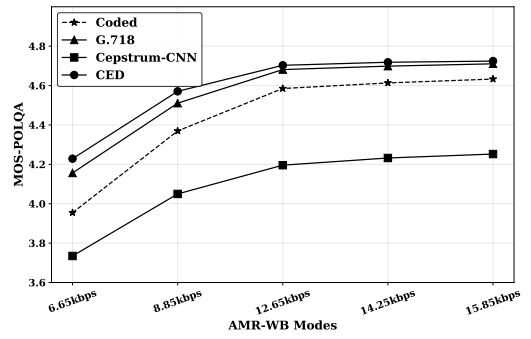


Fig. 6. POLQA scores evaluating the performance of Cepstrum-CNN, CED and G.718 using the TIMIT test set

signal in terms of energy, but at the price of amplifying the coded artefacts.

- In addition, at high frequencies, our proposed CED model has better capability of restoring the harmonic structure lost due to high quantization noise compared to the Cepstrum-CNN.

In order to test the model on completely unseen data, we performed a cross-database validation test. Fig 6 shows the performance of all the models using the test set of TIMIT database [24]. The behavior is similar to the test set of the NTT database, thus validating that the proposed model works well on completely unknown data.

5. CONCLUSION

A convolutional encoder-decoder (CED) based post-filter that estimates a real-valued mask per time frequency bin is proposed to enhance the quality of the coded speech. It is shown that modified-signal approximation is necessary to train a generalized model that works well at higher bitrates despite being trained on the lowest one. Based on POLQA and MUSHRA scores, it was confirmed that real-valued mask based post-filter based on data driven approach that makes no assumption about signal or noise characteristics can successfully enhance the quality of the coded speech. The benefits are higher at low bitrates and as the bitrate increases, the benefits observed are smaller. Cross-database testing also confirmed the robustness of our proposed CED model.

6. REFERENCES

- [1] 3GPP, “Speech codec speech processing functions; Adaptive Multi-Rate - Wideband (AMR-WB) speech codec; Transcoding functions,” 3rd Generation Partnership Project (3GPP), TS 26.190, 12 2009. [Online]. Available: <http://www.3gpp.org/ftp/Specs/html-info/26190.htm>
- [2] —, “TS 26.445, EVS Codec Detailed Algorithmic Description; 3GPP Technical Specification (Release 12),” 3rd Generation Partnership Project (3GPP), TS 26.445, 12 2014. [Online]. Available: <http://www.3gpp.org/ftp/Specs/html-info/26445.htm>
- [3] Juin-Hwey Chen and A. Gersho, “Adaptive postfiltering for quality enhancement of coded speech,” *IEEE Transactions on Speech and Audio Processing*, vol. 3, no. 1, pp. 59–71, Jan 1995.
- [4] ITU-T Recommendation G.718, “Frame error robust narrowband and wideband embedded variable bit-rate coding of speech and audio from 8–32 kbit/s,” 2008.
- [5] M. Dietz, M. Multus, V. Eksler, V. Malenovsky, E. Norvell, H. Poblath, L. Miao, Z. Wang, L. Laaksonen, A. Vasilache, Y. Kamamoto, K. Kikuri, S. Ragot, J. Faure, H. Ehara, V. Rajendran, V. Atti, H. Sung, E. Oh, H. Yuan, and C. Zhu, “Overview of the evs codec architecture,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2015, pp. 5698–5702.
- [6] T. Vaillancourt, R. Salami, and M. Jelínek, “New post-processing techniques for low bit rate celp codecs,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2015, pp. 5908–5912.
- [7] S. Das and T. Bäckström, “Postfiltering using log-magnitude spectrum for speech and audio coding,” in *Proc. Interspeech 2018*, 2018, pp. 3543–3547. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2018-1027>
- [8] K. Han, Y. Wang, D. Wang, W. S. Woods, I. Merks, and T. Zhang, “Learning spectral mapping for speech dereverberation and denoising,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 6, pp. 982–992, June 2015.
- [9] Y. Zhao, D. Wang, I. Merks, and T. Zhang, “Dnn-based enhancement of noisy and reverberant speech,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 6525–6529.
- [10] Y. Wang, A. Narayanan, and D. Wang, “On training targets for supervised speech separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, pp. 1849–1858, 2014.
- [11] F. Weninger, J. R. Hershey, J. Le Roux, and B. Schuller, “Discriminatively trained recurrent neural networks for single-channel speech separation,” in *2014 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, Dec 2014, pp. 577–581.
- [12] Z. Zhao, S. Elshamy, H. Liu, and T. Fingscheidt, “A cnn post-processor to enhance coded speech,” in *2018 16th International Workshop on Acoustic Signal Enhancement (IWAENC)*, Sep. 2018, pp. 406–410.
- [13] Z. Zhao, H. Liu, and T. Fingscheidt, “Convolutional neural networks to enhance coded speech,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 4, pp. 663–678, April 2019.
- [14] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. W. Senior, and K. Kavukcuoglu, “Wavenet: A generative model for raw audio,” 2016. [Online]. Available: <http://arxiv.org/abs/1609.03499>
- [15] J. Valin and J. Skoglund, “Lpcnet: Improving neural speech synthesis through linear prediction,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2019, pp. 5891–5895.
- [16] J. Skoglund and J.-M. Valin, “Improving opus low bit rate quality with neural speech synthesis,” in *arXiv preprint arXiv:1905.04628*, 2019.
- [17] ITU-T Recommendation G.114, “One-way transmission time,” 2003.
- [18] *Perceptual objective listening quality assessment (POLQA)*, ITU-T Recommendation P.863, 2011. [Online]. Available: <http://www.itu.int/rec/T-REC-P.863/en>
- [19] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *arXiv preprint arXiv:1412.6980*, 2014.
- [20] ITU-T G.191, “Software tools for speech and audio coding standardization,” 2005.
- [21] ITU-T P.56, “Objective measurement of active speech level,” 2011.
- [22] Recommendation BS.1534, *Method for the subjective assessment of intermediate quality levels of coding systems*, ITU-R, 2003.
- [23] NTT-AT, “Super wideband stereo speech database,” <http://www.ntt-at.com/product/widebandspeech>, accessed: 09.09.2014. [Online]. Available: <http://www.ntt-at.com/product/widebandspeech>
- [24] J. S. Garofolo, L. D. Consortium *et al.*, *TIMIT: acoustic-phonetic continuous speech corpus*. Linguistic Data Consortium, 1993.