

LIP-TO-SPEECH SYNTHESIS IN THE WILD WITH MULTI-TASK LEARNING

Minsu Kim*, Joanna Hong*, Yong Man Ro†

School of Electrical Engineering, KAIST, South Korea

ABSTRACT

Recent studies have shown impressive performance in Lip-to-speech synthesis that aims to reconstruct speech from visual information alone. However, they have been suffering from synthesizing accurate speech in the wild, due to insufficient supervision for guiding the model to infer the correct content. Distinct from the previous methods, in this paper, we develop a powerful Lip2Speech method that can reconstruct speech with correct contents from the input lip movements, even in a wild environment. To this end, we design multi-task learning that guides the model using multimodal supervision, *i.e.* text and audio, to complement the insufficient word representations of acoustic feature reconstruction loss. Thus, the proposed framework brings the advantage of synthesizing speech containing the right content of multiple speakers with unconstrained sentences. We verify the effectiveness of the proposed method using LRS2, LRS3, and LRW datasets.

Index Terms— Lip-to-speech synthesis, Multi-task learning, Multimodal learning, Speech reconstruction, Lip reading

1. INTRODUCTION

With the recent development of Artificial Intelligence (AI) technology, interest in solving problems by connecting AI and humans is increasing to help human life. It is also necessary in human-to-human conversations in everyday life, especially when the importance of virtual meetings and video conferencing is highlighted. Among many problems in human-to-human conversations, the need for technologies recognizing an accurate conversation when voice signals are hardly available has been increasing. This technology is promising since it can help people understand conversation in situations like crowded shopping mall, party with lots of people and loud music, and silent video conference.

There have been ongoing studies in speech reconstruction and speech recognition with visual information only. However, since they solely depend on the visual information (*i.e.*, lip movements) which holds incomplete information about speech [1], these techniques are known to be challenging problems. Especially, video-driven speech reconstruction, also known as Lip-to-speech synthesis (Lip2Speech), has

been shown much lower performance compared to visual speech recognition [2, 3, 4, 5]. Therefore, Lip2Speech is being developed with constrained datasets [6, 7] such that the number of speakers is limited or the sentences following a fixed grammar. In contrast, visual speech recognition [8, 9, 10] achieved significant performance improvements in the wild datasets [11, 12, 13] containing large speaker variations and utterances. The performance gap between the two technologies is because Lip2Speech needs to consider the varying characteristics in speech (*e.g.*, voice, accent, and intonation). Also, in Lip2Speech task, the reconstruction criteria of continuous audio representation is insufficient compared to that of visual speech recognition, while visual speech recognition can be trained with discrete supervision of text, unnecessary to consider the speaker characteristics.

To reduce the performance gaps between Lip2Speech and visual speech recognition, in this paper, we develop a powerful Lip2Speech method that can correctly capture spoken words even in wild environments. To this end, we propose a multi-task learning method that learns to predict text with content supervision and learns to predict acoustic features. Therefore, we can complement the insufficient content guidance in acoustic feature reconstruction criteria through the proposed content supervision. Specifically, we propose two different types of content supervision: feature- and output-level. In the feature-level content supervision, the model is guided to predict aligned text from the input visual representations with Connectionist Temporal Classification (CTC) loss [14] before synthesizing the acoustic features. Therefore, we can bring strong supervision of text into Lip2Speech without losing the alignment of visual inputs and audio outputs. In the output-level content supervision, we adopt an Automatic Speech Recognition (ASR) model to feedback the model to synthesize speech containing the correct words. Along with the proposed content supervision, the model is guided by reconstruction loss which is applied at continuous auditory features, thus forming multi-task learning. Hence, through the proposed framework, it is possible to synthesize high-quality speech by watching only the visual information even in the wild datasets which are composed of utterances from hundreds of speakers and unconstrained sentences.

To validate the effectiveness of the proposed method, we utilize LRS2 [12] and LRS3 [13], the largest sentence-level audio-visual datasets obtained in the wild, and LRW [11],

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. NRF-2022R1A2C2005529). *Equal Contribution. †Corresponding Author.

a word-level dataset. Through comprehensive experiments, we show that the synthesized speech of the proposed method contains correct contents with the lowest Word Error Rates (WERs) compared to the previous state-of-the-art methods.

2. RELATION TO PRIOR WORK

Early works [15, 2] utilized constrained datasets to develop the Lip2Speech models. Later, [16, 5, 4, 17] improved the architecture and training methods, and provided the possibility of unconstrained Lip2Speech. Recently, [18] utilized the LRS3 dataset which has no restriction in both the number of speakers and sentences. However, they failed to successfully measure WER for the LRS3 dataset due to difficulties in synthesizing accurate speech in unconstrained sentence-level dataset. Different from the previous works, in this paper, we focus on synthesizing speech with accurate contents by proposing content supervisions at different levels.

3. METHOD

Let $X = \{x_1, \dots, x_T\} \in \mathbb{R}^{T \times H \times W \times C}$ be an input video containing lip movements, $Y = \{y_1, \dots, y_S\} \in \mathbb{R}^{K \times S}$ be acoustic feature of ground-truth speech constructed with mel-spectrogram, and $U = \{u_1, \dots, u_L\} \in \mathbb{R}^L$ be the ground-truth transcription of the utterance with L tokens. Here, T indicates the frame lengths, H , W , and C are the frame height, width, and channel sizes, respectively, K and S are the mel-spectral channel and the sequence length, respectively, and L represents the length of the transcription. Our main goal is to translate the input lip video X into adequate speech Y without constraint on the number of speakers or sentences (*i.e.*, in the wild). To this end, we guide the model with multi-tasks, text prediction using content supervision and acoustic feature prediction using reconstruction supervision. For the content supervision, we employ two different types of supervision, feature- and output-level. The details of the aforementioned criteria will be described in the following subsections. The overall proposed framework is illustrated in Fig. 1.

3.1. Feature-level content supervision

To synthesize speech by watching the lip movements only, predicting the right spoken words in advance is important. However, most previous works [15, 2, 5, 16, 4] depend on the reconstruction loss (*e.g.*, L1 and L2 loss) between the predicted and ground-truth acoustic representations (*e.g.*, MFCC or mel-spectrogram). Different from the discrete nature of text, meaning that the same words always have the same labels, acoustic representations might contain different values, based upon the different characteristics of the speakers, tones, accents, and so on, for the same words [19]. Thus, using only reconstruction loss for continuous acoustic features for Lip2Speech can be insufficient for guiding the network to infer the right words from the silent talking face video. In or-

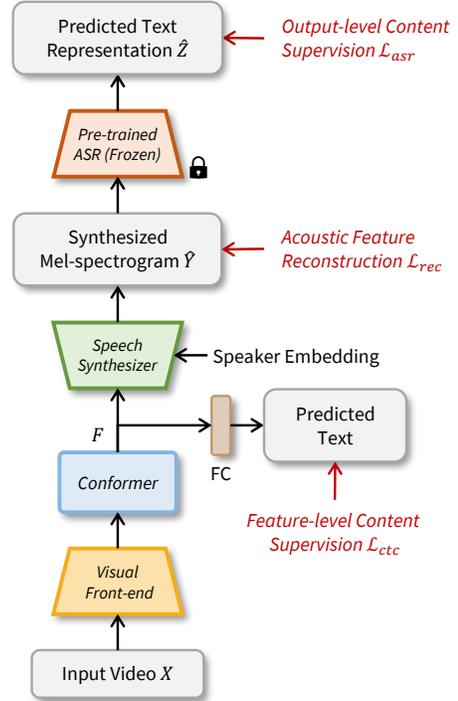


Fig. 1. Overview of the proposed multi-task learning.

der to mitigate this problem, we additionally utilize the discrete text modality with Connectionist Temporal Classification (CTC) loss [14]. Since CTC loss not only guides the network to infer the right words but also produces the aligned representations, it is adequate for the Lip2Speech task which should maintain the visual-audio synchronization.

The frames of input video X is embedded with CNN-based visual front-end Φ_v and their temporal relationships are modeled by conformer Φ_c [20], as follows,

$$F = \Phi_c(\Phi_v(X)). \quad (1)$$

The final encoded visual representations $F = \{f_1, \dots, f_T\} \in \mathbb{R}^{T \times D}$ are employed for both text prediction and speech prediction, where D is the embedding size. The text prediction $p = \text{Softmax}(FW_{ctc} + b_{ctc})$ is guided with the CTC loss defined as $\mathcal{L}_{ctc}(p, U)$, where $W_{ctc} \in \mathbb{R}^{D \times N}$ and $b_{ctc} \in \mathbb{R}^N$ are the weight and bias for text prediction, respectively, and N refers to the number of classes. With the feature-level content supervision, the visual representations F can contain the correct wordings of the input visual speech, which will be eventually translated into output acoustic features.

3.2. Output-level content supervision

As we discussed in the previous section, using CTC loss before synthesizing the speech can explicitly guide the model to learn the correct words. In addition, along with the feature-level content supervision, we can also impose content supervision at the output level. To this end, we propose to bring a

feedback network that guides the model to focus on modelling accurate content even when the text annotation is unavailable. We utilize a pre-trained ASR model to examine whether the synthesized speech contains correct words or not. The content representations \hat{Z} of synthesized speech \hat{Y} are extracted from the ASR model before the last classification layer. Then, the content representations \hat{Z} are guided to resemble with the ground-truth content representations Z which is extracted from the ground-truth speech Y as follows,

$$Z = \text{ASR}(Y), \quad \hat{Z} = \text{ASR}(\hat{Y}), \quad (2)$$

$$\mathcal{L}_{asr} = \|Z - \hat{Z}\|_2^2. \quad (3)$$

Therefore, the model can focus more on modelling accurate content in the output with the help of the feedback model. Moreover, since the output-level content supervision does not require any text annotation, we can employ it even if the text annotations are unavailable.

3.3. Lip-to-Speech synthesis in the wild

Different from the previous work, [18], we put the speaker embedding [16], which for providing the speaker characteristics, right before the speech synthesizer so that the front models (*i.e.*, visual front-end and conformer) can focus more on modelling speech content. Moreover, to properly embed the speaker characteristics into the final output speech, we use a deeper speech synthesizer, composed of 1D CNN layers.

From the visual features F , acoustic features are generated by Speech Synthesizer Ψ . The generation is guided with the reconstruction loss as follows,

$$\mathcal{L}_{rec} = \|\Psi(F) - Y\|_1. \quad (4)$$

By weighted summing the three loss functions defined, the objective function for the proposed multi-task learning can be formulated as follows, $\mathcal{L}_{tot} = \lambda_{ctc}\mathcal{L}_{ctc} + \lambda_{asr}\mathcal{L}_{asr} + \lambda_{rec}\mathcal{L}_{rec}$, where λ are the balancing weights. By guiding the network with multimodal supervision, the synthesized speech $\hat{Y} = \Psi(F)$ can contain accurate contents with the help of content supervision. For stable training, we use 0 for λ_{asr} for early epochs and turn on the output-level content supervision when the reconstruction loss falls enough.

4. EXPERIMENTS

4.1. Datasets

LRS2 [12] dataset is utilized to validate the performance of the proposed method in the wild environment without constraints of the number of speakers and sentences. It has about 142,000 utterances including pre-train and train sets. We utilize both sets for training (about 223 hours), and test the model on a test set containing 1,243 utterances.

LRS3 [13] dataset is another large-scale audio-visual corpus dataset. We follow the unseen data splits of [18]. For training,

about 131,000 utterances are utilized (about 296 hours), and 1,308 utterances are used for testing.

LRW [11] dataset is a word-level English audio-visual corpus. It contains 500 word classes with a maximum of 1,000 training videos each. The dataset is collected from BBC news program, thus containing large speaker, pose, and illumination variations. Compared to the previously mentioned sentence-level datasets in the wild, it is constrained to 500 words. The total training data length is about 157 hours.

4.2. Implementation details

4.2.1. Dataset preprocessing

For all datasets, we crop the lip regions, resize the cropped frames into 112×112 , and convert the colors into grayscale as [4]. Audio with a sample rate of 16kHz is converted to mel-spectrogram by using a hop length of 10ms and window length of 40ms. For the data augmentation, we randomly erase the spatial region of the video consistently in all frames, and time-masking [18] is applied to the input video which is beneficial in modelling visual context. For the text tokenizer, we employ sentencepiece [21] for the LRS2 and LRS3 with 4,000 subwords, and word class for the LRW dataset.

4.2.2. Architectural details

For the visual front-end, ResNet18 whose first layer is modified with a 3D convolution is utilized [22]. For the conformer, we use 12 layers, 8 heads, and convolution kernel size of 31, for the LRS2 and LRS3. For the LRW, we use 6 layered conformer. For the speech synthesizer, we use 3 layered 1D convolutions with 256, 128, and 320 hidden dimensions and a kernel size of 7, and the speaker embedding [16, 18] is obtained by encoding random short audio clip of the same speaker with another 3 layered 1D convolutions with 128, 256, and 256 hidden dimensions and a kernel size of 7. For the LRW, we use random 0.2 sec of audio clip and 0.5 sec for both LRS2 and LRS3. The speaker embedding setting is different from the previous work [18] that utilized all audio frames to extract the speaker embedding by using a pre-trained speaker verification model on large-scale extra datasets. For each LRS2 and LRS3 dataset, the pre-trained ASR model is trained using CTC loss and is composed of 6 conformer encoders with 4 attention heads. For the LRW dataset, we utilize temporal average pooling and cross entropy loss. The pre-trained ASR models keep being frozen during training the Lip2Speech model.

4.2.3. Training details

For training, we use an initial learning rate of 0.0001, batch size of 16, and AdamW [23] optimizer. For the LRS2 and LRS3 datasets, mel-spectrogram of 200 frames are generated by randomly selecting corresponding visual representations from F , to save computing memory. λ_{ctc} , λ_{asr} , and λ_{rec} are empirically set to 1, 1, and 100, respectively.

Table 1. Ablation study on LRS2 dataset. ‘C.S.’ is an abbreviation for content supervision

Method	STOI	ESTOI	PESQ	WER(%)
Baseline [18]	0.491	0.291	1.34	93.91
+ Feature-level C.S.	0.513	0.312	1.34	76.40
+ Conv1D	0.518	0.320	1.34	75.32
+ Speaker Embedding	0.519	0.329	1.36	68.80
+ Output-level C.S. (Full)	0.526	0.341	1.36	60.54
- Feature-level C.S.	0.496	0.299	1.34	84.27

Table 2. Performance comparisons on LRS2 dataset

Method	STOI	ESTOI	PESQ	WER(%)
VCA-GAN [4]	0.407	0.134	1.24	109.01
SVTS [18]	0.491	0.291	1.34	93.91
Proposed Method	0.526	0.341	1.36	60.54

4.2.4. Evaluation metrics

For evaluation, we use Short-Time Objective Intelligibility (STOI) [24] and Extended STOI (ESTOI) [25] to measure the intelligibility of generated speech, Perceptual Evaluation of Speech Quality (PESQ) [26] to measure the perceptual quality of generated speech, and Word Error Rate (WER) to measure how the generated speech containing the correct contents. Note that we focus on WER metric to examine whether the proposed method is effective in modelling speech content. The quantitative results are obtained by converting into the waveform with Griffin-Lim algorithm [27]. Note that the ASR model used to train the Lip2Speech model and that used to measure quantitative results are different.

4.3. Ablation study

We examine the effectiveness of the proposed method by adding each component to the baseline, SVTS [18]. We use the LRS2 dataset to validate the efficacy of the proposed framework in the wild environment. Table 1 shows the ablation results. The WER of baseline model is 93.91%, meaning that the baseline model can hardly generate speech with correct content. By using the feature-level content supervision (+ Feature-level C.S), we improve the WER by 17.51% which is a large improvement compared to the baseline model. Moreover, by changing the speech synthesizer with a deeper architecture (+ Conv1D) and the injection location of speaker embedding (+ Speaker Embedding), the performance is improved to 68.80% WER. Finally, with the proposed output-level content supervision (+ Output-level C.S) which is the final proposed method (*i.e.*, Full), we achieve 60.54% WER, and the proposed model outperforms the baseline model by 33.37% WER. Finally, in the case of that text annotation is not available (- Feature-level C.S), we achieve 84.27% WER with the output-level content supervision, which improves the WER by about 10% from the baseline.

The final performance (60.5% WER) is significant since

Table 3. Performance comparisons on LRS3 dataset

Method	STOI	ESTOI	PESQ	WER(%)
VCA-GAN [4]	0.474	0.207	1.23	96.63
SVTS [18]	0.507	0.271	1.25	79.83
Proposed Method	0.497	0.268	1.31	66.78

Table 4. Performance comparisons on LRW dataset.

Method	STOI	ESTOI	PESQ	WER(%)
Lip2Wav [16]	0.543	0.344	1.20	34.20 [‡]
VCA-GAN [4]	0.565	0.364	1.34	32.07
End-to-end GAN [5]	0.552	0.330	1.33	41.31
SVTS [18]	0.649	0.483	1.49	12.53
Proposed Method	0.642	0.476	1.56	13.86

a well-known visual speech recognition with CTC architecture [8] achieves about 65% WER on LRS2 dataset. This means that our proposed Lip2Speech model achieves comparable performance with the visual speech recognition tasks.

4.4. Comparisons with previous methods

We compare the performances with the previous methods by using the large-scale audio-visual datasets, LRS2 and LRS3, composed of various utterances from diverse speakers. Since there is no prior work that tried to measure the WER on such wild sentence-level datasets, we train VCA-GAN [4] and SVTS [18] on the LRS2 and measure the WER, shown in Table 2. Compared to the previous methods, VCA-GAN and SVTS, the results confirm that the proposed method synthesizes the speech with accurate wordings by achieving the best WER and PESQ. Moreover, in Table 3, we obtain consistent results for the LRS3 dataset compared to the previous experiments with the large gap in the WER performances.

In addition, we verify that our proposed model achieves comparable performances with the other popular methods using word-level dataset, LRW. As indicated in Table 4, the proposed method achieves comparable performances with the state-of-the-art methods showing that multi-task learning can be also applied to the word-level dataset.

5. CONCLUSION

In this paper, we design a powerful Lip2Speech framework that works in the wild. To this end, we propose multi-task content learning: feature- and output-level content supervisions, with the acoustic feature reconstruction. The extensive experimental results prove that the proposed learning framework effectively translates the lip movements into speech audio with the accurate content, even in the wild environments.

[‡]Reported using Google API.

6. REFERENCES

- [1] Minsu Kim, Joanna Hong, Se Jin Park, and Yong Man Ro, “Cromm-vs: Cross-modal memory augmented visual speech recognition,” *IEEE Transactions on Multimedia*, vol. 24, pp. 4342–4355, 2022.
- [2] Konstantinos Vougioukas, Pingchuan Ma, Stavros Petridis, and Maja Pantic, “Video-driven speech reconstruction using generative adversarial networks,” *arXiv preprint arXiv:1906.06301*, 2019.
- [3] Daniel Michelsanti, Olga Slizovskaia, Gloria Haro, Emilia Gómez, Zheng-Hua Tan, and Jesper Jensen, “Vocoder-based speech synthesis from silent videos,” *arXiv preprint arXiv:2004.02541*, 2020.
- [4] Minsu Kim, Joanna Hong, and Yong Man Ro, “Lip to speech synthesis with visual context attentional gan,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [5] Rodrigo Mira, Konstantinos Vougioukas, Pingchuan Ma, Stavros Petridis, Björn W Schuller, and Maja Pantic, “End-to-end video-to-speech synthesis using generative adversarial networks,” *IEEE Transactions on Cybernetics*, 2022.
- [6] Martin Cooke, Jon Barker, Stuart Cunningham, and Xu Shao, “An audio-visual corpus for speech perception and automatic speech recognition,” *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421–2424, 2006.
- [7] Naomi Harte and Eoin Gillen, “Tcd-timit: An audio-visual corpus of continuous speech,” *IEEE Transactions on Multimedia*, vol. 17, no. 5, pp. 603–615, 2015.
- [8] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman, “Deep audio-visual speech recognition,” *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [9] Pingchuan Ma, Stavros Petridis, and Maja Pantic, “End-to-end audio-visual speech recognition with conformers,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7613–7617.
- [10] Minsu Kim, Jeong Hun Yeo, and Yong Man Ro, “Distinguishing homophones using multi-head visual-audio memory for lip reading,” in *Proceedings of the 36th AAAI Conference on Artificial Intelligence, Vancouver, BC, Canada, 2022*, vol. 22.
- [11] Joon Son Chung and Andrew Zisserman, “Lip reading in the wild,” in *Asian Conference on Computer Vision*. Springer, 2016, pp. 87–103.
- [12] Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman, “Lip reading sentences in the wild,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 3444–3453.
- [13] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman, “Lrs3-ted: a large-scale dataset for visual speech recognition,” *arXiv preprint arXiv:1809.00496*, 2018.
- [14] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [15] Hassan Akbari, Himani Arora, Liangliang Cao, and Nima Mesgarani, “Lip2audspec: Speech reconstruction from silent lip movements video,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 2516–2520.
- [16] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar, “Learning individual speaking styles for accurate lip to speech synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13796–13805.
- [17] Joanna Hong, Minsu Kim, Se Jin Park, and Yong Man Ro, “Speech reconstruction with reminiscent sound via visual voice memory,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3654–3667, 2021.
- [18] Rodrigo Mira, Alexandros Haliassos, Stavros Petridis, Björn W Schuller, and Maja Pantic, “Svts: Scalable video-to-speech synthesis,” *arXiv preprint arXiv:2205.02058*, 2022.
- [19] Aaron Van Den Oord, Oriol Vinyals, et al., “Neural discrete representation learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [20] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.
- [21] Taku Kudo and John Richardson, “Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” *arXiv preprint arXiv:1808.06226*, 2018.
- [22] Stavros Petridis, Themos Stafylakis, Pingchuan Ma, Feipeng Cai, Georgios Tzimiropoulos, and Maja Pantic, “End-to-end audiovisual speech recognition,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6548–6552.
- [23] Ilya Loshchilov and Frank Hutter, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [24] Cees H Taal, Richard C Hendriks, Richard Heusdens, and Jesper Jensen, “A short-time objective intelligibility measure for time-frequency weighted noisy speech,” in *2010 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2010, pp. 4214–4217.
- [25] Jesper Jensen and Cees H Taal, “An algorithm for predicting the intelligibility of speech masked by modulated noise maskers,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 11, pp. 2009–2022, 2016.
- [26] Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra, “Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs,” in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No. 01CH37221)*. IEEE, 2001, vol. 2, pp. 749–752.
- [27] Daniel Griffin and Jae Lim, “Signal estimation from modified short-time fourier transform,” *IEEE Transactions on acoustics, speech, and signal processing*, vol. 32, no. 2, pp. 236–243, 1984.