# Deep Learning for Joint Acoustic Echo and Acoustic Howling Suppression in Hybrid Meetings

*Hao Zhang, Meng Yu, Dong Yu*

Tencent AI Lab, Bellevue, WA, USA

{aaronhzhang, raymondmyu, dyu}@global.tencent.com

## Abstract

Hybrid meetings have become increasingly necessary during the post-COVID period and also brought new challenges for solving audio-related problems. In particular, the interplay between acoustic echo and acoustic howling in a hybrid meeting makes the joint suppression of them difficult. This paper proposes a deep learning approach to tackle this problem by formulating a recurrent feedback suppression process as an instantaneous speech separation task using the teacher-forced training strategy. Specifically, a self-attentive recurrent neural network is utilized to extract the target speech from microphone recordings with accessible and learned reference signals, thus suppressing acoustic echo and acoustic howling simultaneously. Different combinations of input signals and loss functions have been investigated for performance improvement. Experimental results demonstrate the effectiveness of the proposed method for suppressing echo and howling jointly in hybrid meetings.

**Index Terms**: hybrid meetings, acoustic echo cancellation, acoustic howling suppression, teacher forcing training

## 1. Introduction

Hybrid meetings, which involve a combination of in-person and remote participants, have become increasingly essential in the post-COVID era [1, 2]. As of 2022, a significant proportion of workplaces (78%) have adopted hybrid work strategies, indicating a growing trend towards hybrid work as the future of work [3]. However, despite the benefits of hybrid meetings, audio-related problems such as acoustic echo and acoustic howling can pose significant challenges and need to be addressed to ensure full-duplex communication.

Acoustic echo refers to the phenomenon where sound originating from a speaker on one end of a communication system is captured by the microphone on the other end and subsequently replayed back to the speaker, creating an unwanted echoing effect [4, 5]. Acoustic howling arises when sound from the speaker's end is captured by the microphone on the same end, leading to a feedback loop that amplifies the sound until it becomes unbearable [6, 7]. Despite having similar underlying mechanisms, acoustic echo and howling are distinct problems, and they can be particularly challenging to address in hybrid meetings where both issues can occur simultaneously. Therefore, it is crucial to have robust and effective algorithms that can address both acoustic echo cancellation (AEC) and acoustic howling suppression (AHS) in a joint manner, taking into account the complex acoustics of the hybrid meeting environment. However, the presence of one problem can affect the estimation and suppression of the other, making it difficult for conventional algorithms to effectively suppress both echo and howling jointly.

Recently, deep learning has emerged as a promising approach for solving the challenges of AEC and AHS due to its ability to model complex nonlinear relationships [8, 9, 10, 11, 12, 13, 14]. In AEC, the problem can be directly formulated as a supervised speech separation problem [15, 8, 16]. However, AHS poses a more complex challenge since it involves the recursively amplification of the playback signal, which makes formulating it as a supervised learning problem non-trivial. To address this challenge, Zhang et al. [14] recently proposed a deep learning based AHS method (Deep AHS) using teacher-forced training strategy, resulting in improved performance when compared to baselines. We believe that recent advances in deep learning based AEC and AHS make it possible to develop effective deep learning methods to address them jointly and solve the full-duplex communication problem in hybrid meetings.

In this study, we tackle the challenges posed by joint AEC and AHS by considering them as an integrated feedback suppression problem and propose a deep learning approach to address it. The recursive feedback suppression process is converted to a speech separation process through teacher forcing training strategy [17, 18], which simplifies the problem formulation and accelerates model training. To accomplish this task, a self-attentive recurrent neural network (SARNN) [19] is utilized to extract target speech from microphone signal with multiple reference signals as additional inputs. Various combinations of inputs are explored to take full use of the accessible reference signals. Given the difficulties in suppressing both forms of feedback jointly, a specific loss function is designed to mitigate leakage introduced due to improperly suppressed feedback, with results demonstrating its efficacy. Experimental results show the effectiveness of the proposed method for joint echo and howling suppression.

The structure of this paper is as follows: Section 2 provides an overview of the audio-related issues in hybrid meeting systems. Section 3 presents the proposed method. The experimental setup is outlined in Section 4, and Section 5 reports the corresponding results. Section 6 concludes the paper.

## 2. Hybrid meetings

### 2.1. Signal model in a hybrid meeting system

For a hybrid meeting system with $J$ devices on the same end and all of them have both a loudspeaker and a microphone turned on, then the total number of acoustic paths in the system will be $J^2$. Take a simplified system with two devices on the same end as an example, as shown in Figure 1 (a). While capturing the target speech $s_i$, the microphone on device $i$ will also record the
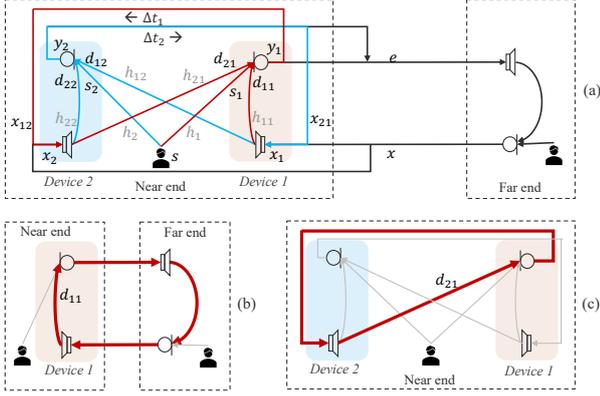
Figure 1: *(a) A simplified hybrid meeting system. (b) and (c) illustrate the two closed acoustic loops related to device 1.*



Figure 2: *Signal flow with teacher-forced learning strategy.*

background noise $n_i$, and playback signals from all devices:

$$y_i = s_i + n_i + \sum_{j=1}^{2} d_{ji} = s_i + n_i + \sum_{j=1}^{2}(x_j * h_{ji}) \quad (1)$$

where $x_j$ is the loudspeaker signal on device $j$, and $d_{ji}$ is the signal picked up by microphone $i$ from loudspeaker $j$ through the acoustic path $h_{ji}$. Among these playback signals, $d_{ii}$ is the playback from device $i$'s own loudspeaker to its microphone, which is known as acoustic echo. Compared to $d_{ji}$ ($j \neq i$), acoustic echo ($d_{ii}$) is relatively easier to suppress since each device usually only has access to its own loudspeaker signal $x_i$, which can be used as a reference signal during the attenuation of $d_{ii}$.

Challenges arise when speakers on the far end and near end talk simultaneously. Considering that each device cannot distinguish whether other devices are exposed in the same space or not, it treats all other devices as far end and sends its processed signal to them. The loudspeaker signal $x_i$ will then be a combination of the far-end signal $x$ and the processed signals sent to device $i$ from device $j$ (denoted as $x_{ji}, (j \neq i)$):

$$x_i = x + x_{ji}, j \neq i \quad (2)$$

If feedback suppression module on each device works properly, the resulting processed signal, $x_{ji}$, should resemble a delayed, scaled, and reverberant version of the near end speech $s$. From the perspective of signal sources, microphone signal given in (1) can be rewritten as:

$$y_i = s_i + n_i + \sum_{j=1}^{2} d_{ji}^x + \sum_{j=1}^{2} d_{ji}^s \quad (3)$$

where $d_{ji}^x$ and $d_{ji}^s$ represent the playback components originated from $x$ and $s$, respectively. It is more challenging to suppress $d_{ji}^s$ because it comes from the same source as that of the target speech $s_i$, and reducing it could distort the target signal.

### 2.2. Joint acoustic echo and acoustic howling suppression

Let us focus on device 1 to analysis the audio-related problem in a hybrid meeting system. There are two closed acoustic loops (CAL) per device in the system, shown in Figure 1 (b) and (c), that can cause acoustic howling. The first CAL, due to acoustic echo, is easier to handle since acoustic echo occurs once per
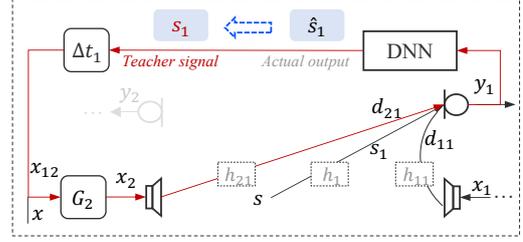
transmission and is handled on both ends. The second CAL is more challenging due to two reasons: 1) Device 1 lacks access to the reference signal that causes feedback $d_{21}$. 2) The two devices involved in this CAL are exposed in the same space.

Without any processing, the microphone signal will be played out through loudspeaker and repeatedly re-enter the pickup. The microphone signal $y_1$ at time index $t$ can then be represented as:

$$y_1(t) = s_1(t) + n_1(t) + d_{11}(t) + \quad (4)$$
$$NL\left[(y_1(t - \Delta t_1) + x(t)) \cdot G_2\right] * h_{21}(t)$$

where $\Delta t_1$ denotes the system delay from device 1 to device 2, $G_2$ is the gain of amplifier on device 2, and $NL(\cdot)$ the non-linear function of loudspeaker. Playback $d_{11}(t)$ is the acoustic echo. The recursive relationship between $y_1(t)$ and $y_1(t - \Delta t)$ causes re-amplifying of playback signal and leads to an annoying, high-pitched sound, which is known as acoustic howling.

In hybrid meetings, achieving full-duplex communication requires addressing both AEC and AHS simultaneously. Nonetheless, the presence of either issue can hinder the accurate detection and elimination of the other, resulting in a shortage of effective solutions.

## 3. Proposed method

### 3.1. Problem formulation

To address the recursive nature of howling, a deep neural network (DNN) module needs to be integrated into the closed acoustic loop and trained recursively. However, this is not practical due to its high computational cost. Alternatively, teacher forcing training strategy can be used to formulate the joint AEC and AHS task as a general feedback suppression problem, as is detailed in Figure 2. This is based on the assumption that the model, once properly trained, can attenuate all feedback signals ($d_{11}$ and $d_{21}$) and transmit only the target speech $s_1$. Through teacher-forced learning, the actual output $\hat{s}_1$ is replaced with the teacher signal $s_1$ during model training. As a result, rather than generated recursively, the microphone signal (4) is simplified to a mixture of target signal, background noise, acoustic echo, and an one-time playback signal determined by $s_1$:

$$y_1(t) = s_1(t) + n_1(t) + d_{11}(t) + \quad (5)$$
$$NL[(s_1(t - \Delta t_1) + x(t)) \cdot G_2] * h_{21}(t)$$

And the overall problem can thus be formulated as a speech separation problem during model training where the task is to separate target signal from the microphone recording with accessible loudspeaker signals ($x_1$, and/or $x, x_{21}$) as references.

### 3.2. Inputs and reference signals

Appropriate reference signals, which enables accurate estimation of the playback signals, are crucial for AEC and AHS algorithms. The reference signal for device 1 is a mixture of two signals, as shown in Figure 1 and (2). The most direct approach is to use the integrated signal $x_1$ as a reference for suppressing the two feedback signals $d_{11}$ and $d_{21}$ in $y_1$. However, this may be less effective for suppressing $d_{21}$. Known from equation (3) that the playback signals share common components originating from different sound sources. Depending on the design of the audio system, we could also have access to $x$ and $x_{21}$ in addition to the integrated loudspeaker signal $x_1$. Using separated loudspeaker signals ($x$ and $x_{21}$) as references could make the suppression of both feedbacks more efficient.

Besides these accessible reference signals obtained directly from device, we have also designed the network to estimate some intermediate outputs from the inputs and use them as non-linear reference signals to further improve feedback cancellation performance [9, 20].

### 3.3. Network structure

Network of the proposed method is given in Figure 3. It takes the microphone signal and one or two reference signals (represented as $r_1$ and $r_2$) as inputs. The input signals, sampled at 16 kHz, are transformed into the frequency domain using a 512-point short-time Fourier transform (STFT) with a frame size of 32 ms and frame shift of 16 ms. The resulting frequency domain inputs are labeled as $Y$, $R_1$, and $R_2$, respectively.

To extract more information from inputs and facilitate the suppression of playback signals, we follow [19, 14] and design the input feature as a concatenation of the normalized log-power spectra (LPS), correlation matrix across time frames and frequency bins, and channel covariance of input signals. These features are concatenated and then passed through a linear layer for feature fusion, followed by a gated recurrent unit (GRU) layer with 257 hidden units and three 1D convolution layers to estimate three complex-valued filters. The filters are then applied to the inputs through deep filtering [21] to obtain the corresponding intermediate signals, $\tilde{Y}$, $\tilde{R}_1$, and $\tilde{R}_2$. These signals serve as additional nonlinear reference signals and their LPS are then concatenated with the original fused feature, and another linear layer is used for feature fusion.

Next, an SARNN module is used to estimate a four-channel enhancement filter, which is then applied on the microphone signal and the three learned reference signals to obtain the enhanced target signal $\hat{S}_1$. Finally, an inverse STFT (iSTFT) is used to obtain the waveform $\hat{s_1}$. More details regarding the feature design and network structure can be found in [19].

### 3.4. Loss functions

In the initial stage of this study, a combination of time-domain scale-invariance signal-to-distortion ratio (SI-SDR) [22] loss and frequency-domain mean absolute error (MAE) of spectrum magnitude is used as loss function for model training:

$$Loss_1 = -\text{SI-SDR}(\hat{s}, s) + \lambda\text{MAE}(|\hat{S}|, |S|) \qquad (6)$$

Given that the feedback signals have a strong correlation with the target signal, suppressing them could be difficult. To further suppress the leakage introduced due to improperly attenuated playback signals, we propose to include a correlation loss:

$$Loss_{corr} = [1 - corr(\hat{s}_1, s_1)] + corr(\hat{s}_1 - s_1, d_*) \qquad (7)$$
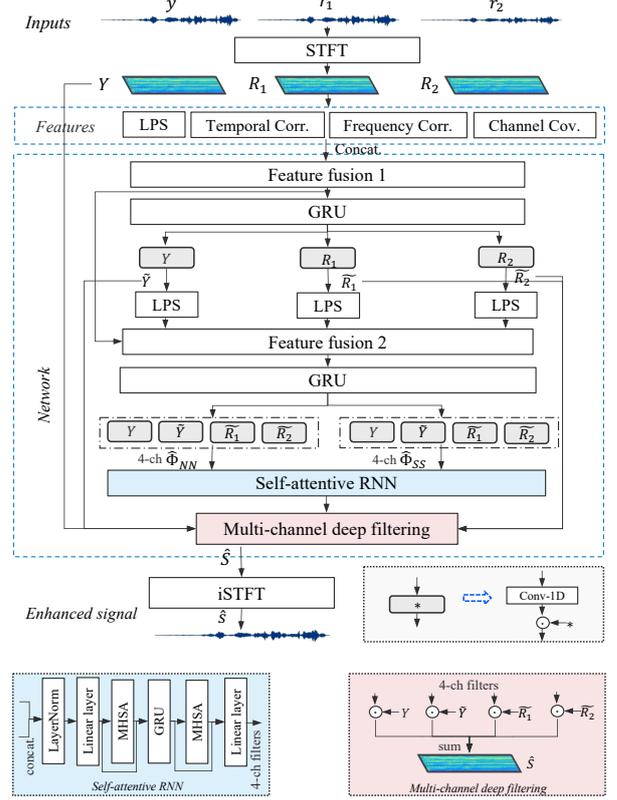


Figure 3: *Architecture of the DNN module for joint acoustic echo and howling suppression. Where the block in gray denotes a combination of 1D convolution layer and deep filtering. A "Conv-1D" outputs a complex-valued ratio filter, which is then applied upon signal * through deep filtering, denoted as ⊙.*

The correlation loss is composed of two terms. The first term evaluates the similarity between the estimated and target signals, while the second term measures the similarity between a playback signal $d_*$ and the residual signal in the estimated target. The modified loss function we used for model training is:

$$Loss_2 = Loss_1 + \beta Loss_{corr} \qquad (8)$$

To ensure balance among different losses, we set $\lambda$ and $\beta$ to 10000 and 10, respectively, in our implementation.

## 4. Experimental setup

### 4.1. Data preparation

The AISHELL-2 [23] and INTERSPEECH 2021 AEC Challenge [24] datasets are used for carrying out experiments. A total number of 10,000 room impulse response (RIR) sets are generated using the image method [25], which incorporate random room characteristics and reverberation times (RT60) range of 0 to 0.6 seconds. Each RIR set consists of 6 RIRs, as shown in Figure 1. During data generation, a randomly chosen RIR set is utilized to create near-end speech signals and the corresponding playback signals. System delay is defined as a random value within the range of $[0.1, 0.3]$ second and microphone nonlinear distortions are simulated using a saturation type of nonlinearity with hard clipping and Sigmoidal function [26, 15, 8]. The microphone signal is generated as a mixture with a randomly cho-

Table 1: *Explorations regarding inputs/reference signals.*

| $Net_{2-ch}, Loss_1$ | SI-SDR (dB) | | | PESQ | | |
|---|---|---|---|---|---|---|
| SFR (dB) | -10 | -5 | 0 | -10 | -5 | 10 |
| Unprocessed | -9.49 | -4.47 | 0.54 | 1.35 | 1.64 | 2.05 |
| $[y_1, x_1]$ | 4.59 | 7.78 | 10.72 | 2.26 | 2.59 | 2.88 |
| $[y_1, x]$ | 3.24 | 6.40 | 9.28 | 2.06 | 2.38 | 2.71 |
| $[y_1, x_{21}]$ | 4.67 | 7.69 | 10.43 | 2.55 | 2.88 | 3.12 |
| $[y_1, x, x_{21}]$ | 5.25 | 8.21 | 10.98 | 2.58 | 2.89 | 3.11 |
| $[y_1, x_{21}, x]$ | **5.31** | **8.53** | **11.42** | **2.69** | **2.99** | **3.23** |

Table 2: *Explorations regarding loss functions.*

| $Net_{2-ch}, [y_1, x_1]$ | | SI-SDR (dB) | | | PESQ | | |
|---|---|---|---|---|---|---|---|
| SFR (dB) | | -10 | -5 | 0 | -10 | -5 | 10 |
| Unprocessed | | -9.49 | -4.47 | 0.54 | 1.35 | 1.64 | 2.05 |
| $Loss_1$ | | 4.59 | 7.78 | 10.72 | 2.26 | 2.59 | **2.88** |
| $Loss_2$ | $d_{21}^s$ | 3.95 | 6.58 | 8.52 | 2.13 | 2.43 | 2.67 |
| | $\tilde{d}_{21}$ | **4.96** | **7.92** | **10.83** | **2.27** | **2.60** | **2.88** |
| | $d_{21}^s + d_{11}^s$ | 4.32 | 6.98 | 9.03 | 2.26 | 2.56 | 2.79 |
| | $d_{21} + d_{11}$ | 4.74 | 7.65 | 10.32 | 2.27 | 2.58 | 2.86 |

Table 3: *Proposed method for feedback suppression.*

| Input: $[y_1, x_{21}, x]$ | SI-SDR | | | PESQ | | |
|---|---|---|---|---|---|---|
| SFR (dB) | -10 | -5 | 0 | -10 | -5 | 0 |
| Unprocessed | -9.49 | -4.47 | 0.54 | 1.35 | 1.64 | 2.05 |
| $Net_{2-ch}, Loss_1$ | 5.31 | 8.53 | 11.42 | 2.69 | 2.99 | 3.23 |
| $Net_{2-ch}, Loss_1$ | 6.11 | 9.09 | 11.82 | **2.73** | **3.07** | **3.31** |
| $Net_{3-ch}, Loss_2$ | **6.47** | **9.24** | **11.87** | **2.73** | 3.03 | 3.28 |

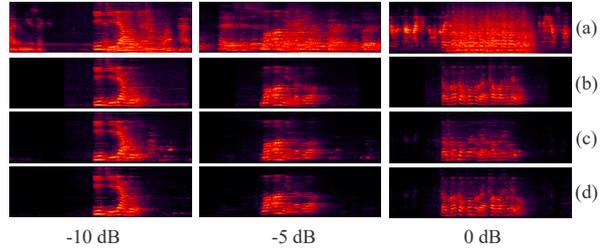

-10 dB          -5 dB          0 dB

Figure 4: *Spectrograms of processed signals obtained under different SFR levels: (a) unprocessed signal, (b) target signal, (c) output of the "initial" model, and (d) output of the best model.*

sen signal-to-feedback ratio (SFR) in the range of $[-20, 5]$ dB, and signal-to-noise (SNR) ratio ranging from -10 dB to 30 dB. We created a total of 10,000, 300, and 500 utterances for training, validation, and testing, respectively. The utterances and RIRs used for generating testing data are different from those used in the training and validation data. The model was trained for 60 epochs using a batch size of 20.

### 4.2. Method evaluation

SI-SDR and perceptual evaluation of speech quality (PESQ) [27] are used as evaluation metrics to show the playback attenuation performance and quality of enhanced speech. A higher value denotes better performance.

## 5. Experimental results

### 5.1. Explorations regarding inputs/reference signals

To investigate the impact of reference signals on model performance, we conducted experiments with different SFRs, and the results are summarized in Table 1. To diminish the influence of model size difference during the comparison, we remove the deep filtering branches related to $\tilde{R}_2$ in the network and use the resulting simplified network, denoted as "$Net_{2-ch}$", for training the inputs with either 2 or 3 channels. This implies that when using an input with three channels, we extract intermediate signals from the first two channels, while the third channel is only used for feature extraction. Using the integrated loudspeaker signal $x_1$ as a reference is the most straightforward way to train the model, and the model trained with $[y_1, x_1]$ is referred to as the "initial" model. Among the models with 2-channel inputs, the "initial" model achieves better SI-SDR in most cases, while its speech quality (PESQ) is not better than that of using $x_{21}$ as the reference signal. Models trained using separated reference signals (3-channel inputs) consistently outperformed models using an integrated reference signal $x_1$. For $Net_{2-ch}$, the order of reference signals $x$ and $x_{21}$ determines from which $\tilde{R}_1$ is extracted. The model trained with $[y_1, x_{21}, x]$ as input achieves the best overall performance.

### 5.2. Explorations regarding loss functions

Table 2 compares the performance of models trained with $Net_{2-ch}$, input $[y_1, x_1]$, and $Loss_2$ calculated using different playback $d_*$. The results show that incorporating the correlation loss does not consistently lead to performance improvement while using $Loss_2$ calculated based on $d_{21}$ yields the best performance and outperforms the model trained using $Loss_1$, especially in terms of SI-SDR. This is because $d_{21}$ is more difficult to suppress due to a lack of direct reference signal, and utilizing it in the calculation of $Loss_2$ helps further attenuating the leakage.

### 5.3. Proposed method for joint AEC and AHS

We combine the findings made through explorations regarding inputs and loss functions and train a model to achieve the best feedback suppression. Specifically, we utilize a 3-channel input $[y_1, x_{21}, x]$, $Loss_2$ with $d_{21}$, and the network illustrated in Figure 3, denoted as "$Net_{3-ch}$", for model training. The performance comparison results presented in Table 3 demonstrate that using "$Net_{3-ch}$" results in better performance than using "$Net_{2-ch}$", and employing the modified loss function, $Loss_2$, could further improve playback attenuating performance. We also provide spectrograms of processed signals obtained using the "initial" model and the best-performing model in Figure 4 to further illustrate the efficacy of our proposed approach for joint acoustic echo and acoustic howling suppression.

## 6. Conclusion

We have proposed a deep learning based method for addressing audio-related problems in hybrid meetings. Our proposed method treats acoustic echo and acoustic howling as an integrated feedback problem and achieves simultaneous AEC and AHS using a teacher-forcing learning strategy. By converting the recursive feedback suppression problem into a speech separation problem, an SARNN model is utilized to extract the target speech from microphone recording with multiple reference signals as additional inputs. The impact of input signals, loss functions on joint AEC and AHS performance has been investigated. Future work includes considering practical issues such as computational complexity and investigating using cascaded network to suppress acoustic echo and howling gradually.

# 7. References

[1] B. Saatçi, R. Rädle, S. Rintel, K. O'Hara, and C. Nyland-sted Klokmose, "Hybrid meetings in the modern workplace: stories of success and failure," in *Collaboration Technologies and Social Computing: 25th International Conference, CRIWG+ CollabTech 2019, Kyoto, Japan, September 4–6, 2019, Proceedings 25.* Springer, 2019, pp. 45–61.

[2] B. Z. Hameed, Y. Tanidir, N. Naik, J. Y.-C. Teoh, M. Shah, M. L. Wroclawski, A. B. Kunjibettu, D. Castellani, S. Ibrahim, R. D. da Silva *et al.*, "Will "hybrid" meetings replace face-to-face meetings post COVID-19 era? perceptions and views from the urological community," *Urology*, vol. 156, pp. 52–57, 2021.

[3] R. Carter, "What is a hybrid meeting? An introduction," https://www.uctoday.com/collaboration/what-is-a-hybrid-meeting-an-introduction/.

[4] J. Benesty, T. Gänsler, D. R. Morgan, M. M. Sondhi, S. L. Gay *et al.*, "Advances in network and acoustic echo cancellation," 2001.

[5] G. Enzner, H. Buchner, A. Favrot, and F. Kuech, "Acoustic echo control," in *Academic press library in signal processing.* Elsevier, 2014, vol. 4, pp. 807–877.

[6] R. V. Waterhouse, "Theory of howlback in reverberant rooms," *The Journal of the Acoustical Society of America*, vol. 37, no. 5, pp. 921–923, 1965.

[7] T. Van Waterschoot and M. Moonen, "Fifty years of acoustic feedback control: State of the art and future challenges," *Proceedings of the IEEE*, vol. 99, no. 2, pp. 288–327, 2010.

[8] H. Zhang and D. Wang, "Deep learning for acoustic echo cancellation in noisy and double-talk scenarios," *Proc. Interspeech 2018*, pp. 3239–3243, 2018.

[9] H. Zhang, S. Kandadai, H. Rao, M. Kim, T. Pruthi, and T. Kristjansson, "Deep adaptive AEC: Hybrid of deep learning and adaptive acoustic echo cancellation," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 2022, pp. 756–760.

[10] H. Zhang and D. Wang, "Deep MCANC: A deep learning approach to multi-channel active noise control," *Neural Networks*, vol. 158, pp. 318–327, 2023.

[11] Z. Chen, Y. Hao, Y. Chen, G. Chen, and L. Ruan, "A neural network-based howling detection method for real-time communication applications," in *2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 2022, pp. 206–210.

[12] H. Gan, G. Luo, Y. Luo, and W. Luo, "Howling noise cancellation in time–frequency domain by deep neural networks," in *Proceedings of Sixth International Congress on Information and Communication Technology.* Springer, 2022, pp. 319–332.

[13] C. Zheng, M. Wang, X. Li, and B. C. Moore, "A deep learning solution to the marginal stability problems of acoustic feedback systems for hearing aids," *The Journal of the Acoustical Society of America*, vol. 152, no. 6, pp. 3616–3634, 2022.

[14] H. Zhang, M. Yu, and D. Yu, "Deep AHS: A deep learning approach to acoustic howling suppression," *arXiv preprint arXiv:2302.09252*, 2023.

[15] C. M. Lee, J. W. Shin, and N. S. Kim, "DNN-based residual echo suppression," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.

[16] H. Zhang, K. Tan, and D. Wang, "Deep learning for joint acoustic echo and noise cancellation with nonlinear distortions." in *Interspeech*, 2019, pp. 4255–4259.

[17] R. J. Williams and D. Zipser, "A learning algorithm for continually running fully recurrent neural networks," *Neural computation*, vol. 1, no. 2, pp. 270–280, 1989.

[18] A. M. Lamb, A. G. ALIAS PARTH GOYAL, Y. Zhang, S. Zhang, A. C. Courville, and Y. Bengio, "Professor forcing: A new algorithm for training recurrent networks," *Advances in neural information processing systems*, vol. 29, 2016.

[19] M. Yu, Y. Xu, C. Zhang, S.-X. Zhang, and D. Yu, "NeuralEcho: A self-attentive recurrent neural network for unified acoustic echo suppression and speech enhancement," *arXiv preprint arXiv:2205.10401*, 2022.

[20] Y. Zhang, M. Yu, H. Zhang, D. Yu, and D. Wang, "KalmanNet: A learnable kalman filter for acoustic echo cancellation," *arXiv preprint arXiv:2301.12363*, 2023.

[21] W. Mack and E. A. Habets, "Deep filtering: Signal extraction and reconstruction using complex time-frequency filters," *IEEE Signal Processing Letters*, vol. 27, pp. 61–65, 2019.

[22] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR–half-baked or well done?" in *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 2019, pp. 626–630.

[23] J. Du, X. Na, X. Liu, and H. Bu, "Aishell-2: Transforming mandarin asr research into industrial scale," *arXiv preprint arXiv:1808.10583*, 2018.

[24] R. Cutler, A. Saabas, T. Pärnamaa, M. Loide, S. Sootla, M. Purin, H. Gamper, S. Braun, K. Sørensen, R. Aichner *et al.*, "INTERSPEECH 2021 acoustic echo cancellation challenge." in *Interspeech*, 2021, pp. 4748–4752.

[25] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *The Journal of the Acoustical Society of America*, vol. 65, no. 4, pp. 943–950, 1979.

[26] A. Birkett and R. Goubran, "Nonlinear loudspeaker compensation for hands free acoustic echo cancellation," *Electronics Letters*, vol. 32, no. 12, pp. 1063–1064, 1996.

[27] 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 *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP).*, vol. 2. IEEE, 2001, pp. 749–752.