# MULTI-TASK SUB-BAND NETWORK FOR DEEP RESIDUAL ECHO SUPPRESSION

*Jiayao Sun[1,2], Dawei Luo[2,*], Zhaoxia Li[2], Jindong Li[2], Yukai Ju[1], Yang Li[2]*

[1] Northwestern Polytechnical University, Xi'an, China
[2] Li Auto, China

## ABSTRACT

This paper introduces the SWANT team's entry to the ICASSP 2023 AEC Challenge. We submit a system that cascades a linear filter with a neural post-filter. Particularly, we adopt sub-band processing to handle full-band signals and shape the network with multi-task learning, where dual signal voice activity detection (DSVAD) and echo estimation are adopted as auxiliary tasks. Moreover, we particularly improve the time frequency convolution module (TFCM) to increase the receptive field using small convolution kernels. Finally, our system has ranked 4th in ICASSP 2023 AEC Challenge Non-personalized track.

***Index Terms***— Acoustic echo cancellation, sub-band processing, multi-task learning

## 1. INTRODUCTION

Neural residual echo suppression [1] has achieved excellent performance in the AEC task. With the great success in previous challenges, the 4th AEC Challenge in ICASSP2023 [2] has particularly focused on more difficult *full-band* signals. To address this challenge with lower complexity, we propose a *sub-band* time frequency domain gated convolutional recurrent neural network (S-TFGCRN) approach. After linear filtering, the full-band signals are first processed to several sub-band signals using the pseudo quadrature mirror filter bank (PQMF) and the sub-band signals are fed separately into an encoder-decoder network to remove echo residuals. Finally, the sub-band signals are merged back into the full-band echo-removed signals. Importantly, our network is optimized under a multi-task learning framework, where dual signal VAD and echo estimation are augmented as auxiliary tasks. To better model the intrinsic relationships between harmonics, we introduce an updated time frequency convolution module (U-TFCM) in the encoder network. The effectiveness of the contributions is validated through an ablation study. According to the official results of the challenge, our model ranks 4th in the non-personalized AEC track.
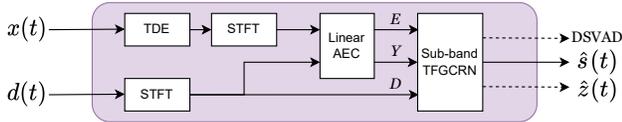


**Fig. 1**. Architecture of our proposed AEC system.

## 2. PROPOSED SYSTEM

### 2.1. Problem formulation

The overall framework of our proposed AEC system is depicted in Fig. 1, where the near-end microphone signal $d(t)$ is a combination

of several signals described as follows.

$$d(t) = s(t) + z(t) + v(t) \tag{1}$$
$$z(t) = x(t) * h(t) \tag{2}$$

where $s(t)$, $x(t)$, $v(t)$ and $z(t)$ denote the near-end speech signal, far-end microphone signal, background noise, and echo signal, respectively, which is generated by the convolution of far-end signal and echo path $h(t)$. The t refers to the time sample index. This task is considered an audio separation task, and the goal is to separate the near-end speech signal $s(t)$ from the near-end microphone signal $d(t)$ and pass it to the far end. The echo $y$ is estimated by adaptive filter NLMS and the signal $e$ is the filter output. $D$, $E$, $Y$, and $Z$ are frequency domain representations of $d$, $e$, $y$, and $z$, respectively. The time delay estimation (TDE) is implemented using Generalized Cross Correlation with PHAse Transform (GCC-PHAT). The dotted line in Fig. 1 indicates that the modules are used as auxiliary tasks and can be removed during model inference.
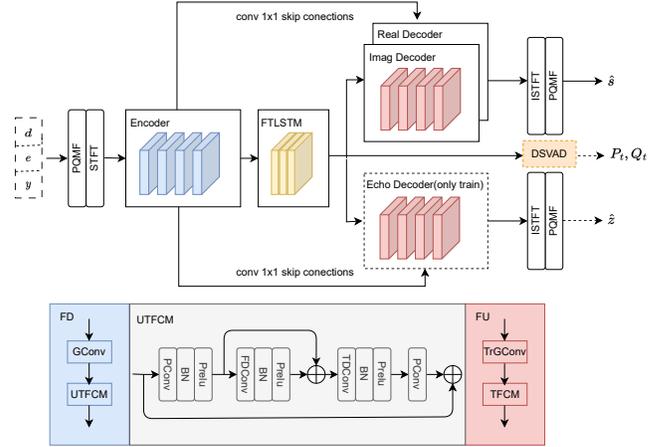


**Fig. 2**. S-TFGCRN architecture.

### 2.2. S-TFGCRN based post-filter

We use the PQMF for signal analysis and synthesis. The sub-band analysis process comprises three steps: FIR analysis, downsampling, and STFT. Likewise, sub-band synthesis involves three stages IS-FTT, upsampling, and FIR synthesis.

As illustrated in Fig. 2, the encoder part contains four frequency downsampling (FD) layers, and every FD layer contains GConv and U-TFCM modules. Different from the original TFCM [3] that mainly focuses on the receptive field of the time dimension, the updated time frequency convolution module (U-TFCM) particularly increases the receptive field of the frequency dimension to learn time-frequency correlation. The skip connection uses a 1x1 convolution. FTLSTM is used as the bottleneck layer, which is proven to be effective in temporal modeling. The structures of GConv and

FTLSTM follow those in [1]. The real/imag decoder part comprises four frequency upsampling (FU) layers, each containing a TrGConv module and a TFCM module. The output of the FU layer is concatenated with the output of the conv1x1. The echo decoder includes 4 FU layers, and its final TrGConv output layer has a channel dimension of 2. Estimating the echo target helps obtain a more accurate estimate of the near-end speech, as we treat the AEC task as an audio separation task. To distinguish between double-talk and single-talk scenarios, we do VAD processing for both near-end and far-end speech, resulting in dual signal VAD (DSVAD) that consists of two VAD modules. The structure of VAD follows that in [1].

### 2.3. Loss Function

Our loss function is based on the echo signal power weighted loss [1], denoted as $\mathcal{L}_{\text{echo-aware}}$. To ensure the accuracy of echo estimation for the echo decoder, $\mathcal{L}_{\text{echo}}$ is defined as

$$\mathcal{L}_{\text{echo}} = \sum_{f,t} |\widehat{Z} - Z|, \tag{3}$$

which is based on the mean absolute error (MAE) in [4]. The binary cross entropy (BCE) is used as the DSVAD loss, formulated as

$$\mathcal{L}_{\text{dtd-nearend}} = \frac{1}{T}\sum_{T} \left(-\bar{P}_t \log(P_t) - (1 - \bar{P}_t) \log(1 - P_t)\right)$$

$$\mathcal{L}_{\text{dtd-farend}} = \frac{1}{T}\sum_{T} \left(-\bar{Q}_t \log(Q_t) - (1 - \bar{Q}_t) \log(1 - Q_t)\right) \tag{4}$$

$$\mathcal{L}_{\text{dtd}} = \mathcal{L}_{\text{dtd-nearend}} + \mathcal{L}_{\text{dtd-farend}}$$

where $t$ and $f$ denote frame and frequency respectively, and $\bar{P}_t, \bar{Q}_t \in \{0, 1\}$ are the near-end and far-end speech activity label respectively, which is based on a short-term energy threshold. $P_t$ and $Q_t$ are the estimated state of near-end and far-end speech, respectively.

The final loss is

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{echo-aware}} + 0.2\mathcal{L}_{\text{mask}} + 0.1\mathcal{L}_{\text{dtd}}$$
$$+ 0.05\mathcal{L}_{\text{echo}} + \mathcal{L}_{\text{asym}} \tag{5}$$

where the definition of $\mathcal{L}_{\text{mask}}$ is the same as that in [1] and $\mathcal{L}_{\text{asym}}$ is consistent with that in [4].

## 3. EXPERIMENTS

### 3.1. Data and experiment setup

In our experiment, the clean speech dataset is from the ICASSP 2022 DNS-challenge [5]. The noise data originates from AudioSet, Freesound, and Demand databases. We also use the far-end single-talk clips officially provided by the AEC challenge as echo signals for training and validation. For RIR, we generate 100,000 pairs using the image method for the training set and 5,000 for the validation set, respectively. The signal-to-noise ratio (SNR) and signal-to-echo ratio (SER) are set to [0, 20]dB and [-10, 15]dB, respectively, for generating near-end microphone signals. SNR is set to [15, 45]dB for far-end signal. We create a small dataset (300 hours) for the ablation study and a large dataset (1400 hours) for final submission model training. In the ablation study, a simulated test set of 1500 clips is used for model evaluation, 400 clips for far-end single talk, 400 clips for near-end single talk, and the rest are double talk.

For the proposed model, the window length and frameshift are 20ms and 10ms respectively. As the computational complexity of the full-band model exceeds the requirements of the challenge, the sub-band approach is taken and the full-band signals are processed into 4 sub-band signals using PQMF. All neural models are trained with the Adam optimizer for 60 epochs with an initial learning rate of 1e-4, and the learning rate is halved if there is no loss decrease on the development set for 2 epochs.

For the output of each layer of encoder/decoder, the number of channels is 80. The kernel size and stride of GConv and TrGConv are set as (2, 3), (1, 2) in the time and frequency axis, respectively. In the U-TFCM module, 4 convolutional layers are adopted with a dilation rate of $\{1, 2, 4, 8\}$ in both the time axis and frequency axis. The stride and kernel size of DConv in U-TFCM are set as (1, 1), (3, 3). The TFCM module is the same as that in [3]. The number of parameters of our submitted model is 3.83M. The RTF of the system is 0.1983 running on Intel(R) Xeon(R) Silver 4214R CPU @ 2.40GHz while the neural post-filter is implemented by ONNX for speed-up.

### 3.2. Results and analysis

We first perform an ablation study on a small data set (300h). Based on the results in Table 1, it is evident that incorporating the DSVAD module results in improvements in all aspects, especially the performance of far-end single-talk suppression. Additionally, integrating the U-TFCM and TFCM modules into the encoder and decoder exhibits extra performance gain, resulting in a further 0.15 WB-PESQ gain for the DT scenario. Moreover, auxiliary training with the echo decoder significantly enhances echo suppression. Finally, we train the whole model using the large dataset (1400h) and process the blind test clips as our submission.

From the blind test set results in Table 2, we can observe that our S-TFGCRN model significantly outperforms the baseline, ranking 4th in the non-personalized track according to the official results.

**Table 1**. Echo suppression performance in the simulated test set. DT: double talk, ST: single-talk, NE: near-end, FE: far-end.

| Model | Para.(M) | DT (WB-PESQ) | ST-NE (WB-PESQ) | ST-FE (ERLE) | Data |
|---|---|---|---|---|---|
| Noisy | | 2.03 | 2.85 | 0 | |
| Sub-band GCRN | 2.34 | 2.84 | 3.29 | 56.58 | |
| + DSVAD | 2.57 | 2.87 | 3.31 | 62.82 | 300h |
| + U-TFCM, TFCM | 3.14 | 3.02 | 3.40 | 63.83 | |
| + Echo Decoder | 3.83 | 3.06 | 3.44 | 66.27 | |
| S-TFGCRN (submitted) | 3.83 | **3.16** | **3.48** | **67.43** | 1400h |

**Table 2**. Model performance on the blind test set.

| Model | Overall MOS | WAcc | Final Score |
|---|---|---|---|
| Baseline | 4.013 | 0.649 | 0.736 |
| S-TFGCRN (submitted) | **4.320** | **0.790** | **0.823** |

## 4. REFERENCES

[1] S. Zhang, Z. Wang, J. Sun, Y. Fu, B. Tian, Q. Fu, and L. Xie, "Multi-task deep residual echo suppression with echo-aware loss," in *ICASSP*. IEEE, 2022.

[2] R. Cutler, A. Saabas, T. Parnamaa, M. Purin, E. Indenbom, N. Ristea, J. Guzhvin, H. Gamper, S. Braun, and R. Aichner, "Icassp 2023 acoustic echo cancellation challenge," 2023.

[3] Y. Ju, S. Zhang, W. Rao, Y. Wang, T. Yu, L. Xie, and S. Shang, "Tea-pse 2.0: Sub-band network for real-time personalized speech enhancement," in *SLT*. IEEE, 2023.

[4] S. Braun and M. Valero, "Task splitting for dnn-based acoustic echo and noise removal," in *IWAENC*. IEEE, 2022.

[5] H. Dubey, V. Gopal, R. Cutler, A. Aazami, S. Matusevych, S. Braun, S. Eskimez, M. Thakker, T. Yoshioka, H. Gamper, et al., "Icassp 2022 deep noise suppression challenge," in *ICASSP*. IEEE, 2022.