

# CabinSep: IR-Augmented Mask-Based MVDR for Real-Time In-Car Speech Separation with Distributed Heterogeneous Arrays

Runduo Han<sup>1</sup>, Yanxin Hu<sup>2</sup>, Yihui Fu<sup>1</sup>, Zihan Zhang<sup>1</sup>, Yukai Jv<sup>1</sup>, Li Chen<sup>2</sup>, Lei Xie<sup>1\*</sup>

<sup>1</sup>Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science, Northwestern Polytechnical University, Xi'an, China

<sup>2</sup> Shanghai ZEEKR Blue New Energy Technology Co., Ltd., China

{rdhan, zhzhzhang, nwpu2016303311} @mail.nwpu.edu.cn, arrowhyx@foxmail.com, felixfuyihui@gmail.com, lichenbook1@163.com, lxie@nwpu.edu.cn

## Abstract

Separating overlapping speech from multiple speakers is crucial for effective human-vehicle interaction. This paper proposes CabinSep, a lightweight neural mask-based minimum variance distortionless response (MVDR) speech separation approach, to reduce speech recognition errors in back-end automatic speech recognition (ASR) models. Our contributions are threefold: First, we utilize channel information to extract spatial features, which improves the estimation of speech and noise masks. Second, we employ MVDR during inference, reducing speech distortion to make it more ASR-friendly. Third, we introduce a data augmentation method combining simulated and real-recorded impulse responses (IRs), improving speaker localization at zone boundaries and further reducing speech recognition errors. With a computational complexity of only 0.4 GMACs, CabinSep achieves a 17.5% relative reduction in speech recognition error rate in a real-recorded dataset compared to the state-of-the-art DualSep model<sup>1</sup>.

**Index Terms:** speech separation, human-car interaction, distortionless constraint, speaker positioning

## 1. Introduction

Speech interaction is crucial for in-car intelligence, with automatic speech recognition (ASR) as a key gateway for human-vehicle interaction. Speech recognition accuracy directly impacts interaction efficiency and user experience [1, 2, 3, 4, 5]. However, when multiple passengers interact with the car simultaneously, overlapping speech becomes inevitable, posing significant challenges to the ASR system [6, 7, 8, 9, 10]. Therefore, it is necessary to apply speech separation systems to separate overlapping speech and distinguish speakers before delivering speech recognition [11].

The neural mask-based minimum variance distortionless response (MVDR) speech separation scheme has been proven effective in improving speech recognition performance by ensuring undistorted speech [12, 13, 14]. However, during the joint training of the neural network and MVDR, the matrix inversion operation in MVDR causes numerical instability, leading to unstable model convergence [15]. Additionally, noise elimination of this scheme is unsatisfactory because MVDR weight selection does not achieve optimal noise suppression [16]. Recent works [16, 17, 18, 19, 20] adopt speech separation schemes that directly use the neural network's output as separated speech, of which [20, 21] approaches achieve remarkable performance in intrusive speech quality metrics including scale-invariant signal-to-noise ratio (SI-SNR)[22]. However, due to the inherent problem of neural networks, nonlinear distortion is intro-

duced in complicated real-world scenarios, negatively affecting ASR performance. Although these advanced models [21, 20] have a relatively small number of parameters, their computational complexity exceeds 1 GMACs, making them challenging to deploy in cars practically. Furthermore, the irregular structures of car cabins make it difficult to simulate training data that matches real-world scenarios. Meanwhile, back-end systems sometimes need to respond based on the speaker's position in in-car speech interactions. As shown in Figure 1(a), accurately locating passengers is challenging when they speak at the boundary of a zone, which previous studies have not adequately addressed.

We propose CabinSep, a plug-and-play streaming speech separation approach for real-world in-car scenarios that improves ASR performance to a great extent. To address nonlinear speech distortion in ASR systems [23, 24], we adopt MVDR during inference to avoid numerical instability and reduce distortion. Furthermore, we introduce a dual-mask estimation mechanism where an auxiliary noise mask complements the speech mask, improving MVDR's noise suppression capabilities. Effectively utilizing spatial information in feature channels aids speech separation tasks, such as the transform-average-concatenate (TAC) module [19, 10]. However, it involves considerable computational complexity. We propose a time skip cascaded TAC module that processes only half of the time frames, reducing computational complexity significantly while maintaining performance. Unlike previous solutions [16, 21], our proposed CabinSep skips the highly accurate direction of arrival (DOA) estimation, which allows for a simpler heterogeneous microphone array design, with each zone in the cabin corresponding to a single-channel microphone, thereby reducing production costs.

To simulate training data that better matches in-car communication scenarios, we propose a data augmentation method that combines simulated impulse responses (IRs) with real-recorded IRs [25]. Notably, this method includes an effective IRs mixing strategy where real-recorded IRs are used for the microphone of the speaker's zone, while simulated IRs are used for other zones. This approach enables CabinSep to accurately localize the zone of the speaker even when the speaker sits at the zone's boundary.

Our proposed CabinSep approach is evaluated by speech recognition accuracy on two open-source pre-trained ASR models, WeNet [26] and SenseVoice [27], which validates the robustness of CabinSep approach across different back-end ASR models architectures. All test data are real-recorded audio from an electric vehicle to better match real-world scenarios. This includes data recorded while the car is stationary or during motion. When the vehicle is moving, the speech is disturbed by complicated background noise, including wind,

\* Corresponding author.

<sup>1</sup>Demos are available at: <https://cabinsep.github.io/cabinsep/>

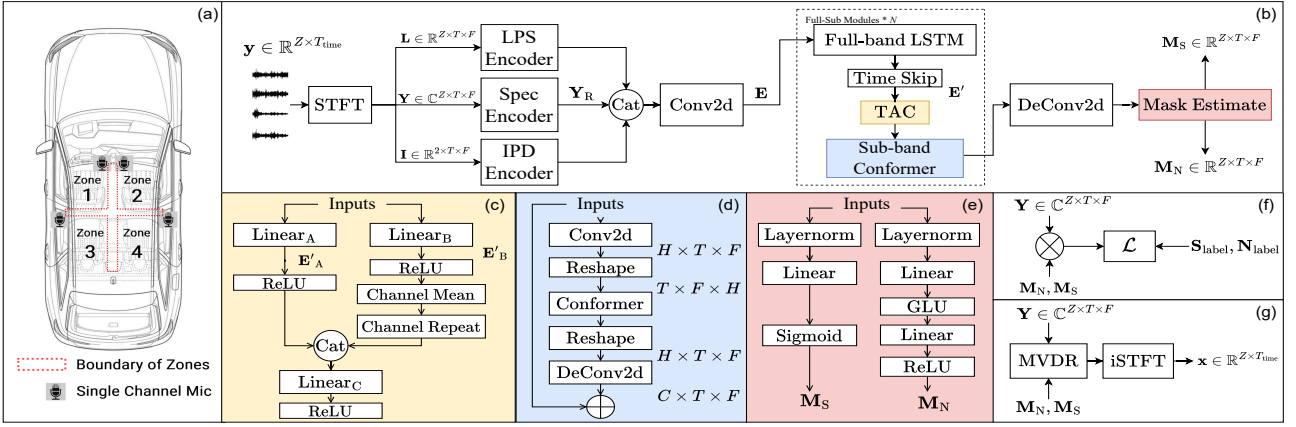


Figure 1: The overall structure of CabinSep. (a) in-car speech separation scenario; (b) model architecture; (c) transform-average-concatenate (TAC) module; (d) sub-band-conformer module; (e) mask estimate module; (f) Training procedure; (g) Inference procedure;  $y$  represents  $Z$  channel audio mixture;  $Y$  represents complex spectrum of  $y$ ;  $L$  represents LPS;  $I$  represents IPD;  $M_S, M_N$  represent estimated speech and noise masks;  $S_{\text{label}}, N_{\text{label}}$  represent speech and noise labels;  $x$  represents separated clean speech;  $\mathcal{L}$  represents the loss function.

wheel, engine, etc [28]. Experimental results show that our proposed CabinSep achieves a 17.5% relative reduction in speech recognition error rate compared to the state-of-the-art (SOTA) approach DualSep [20] for in-car speech separation and recognition, with a computational complexity of only 0.4 GMACs and 0.21 real-time factor (RTF) using the single-core of Qualcomm SA8295P in-car CPU.

## 2. Problem Formulation

As shown in Figure 1(a), we focus on in-car speech separation. We divide the car cabin into  $Z$  zones, with each zone corresponding to a single channel microphone, and at most one person speaks in each zone. Suppose there are  $P$  person speaking in the car ( $P \leq Z$ ), the clean speech corresponding to each zone is set as  $s(z)$ , and the signal recorded by each microphone is  $y_i$ , which can be expressed as:

$$x_i(z) = r_i(z) * s(z), \quad (1)$$

$$y_i = \sum_{z=1}^Z x_i(z) + v_i, \quad (2)$$

where  $z$  represents the zone index in the cabin,  $i \in \{1, \dots, Z\}$  represents the microphone index,  $r$  represents the IRs, and  $v$  is the background noise received by microphone. We aim to estimate the separated clean speech  $x_z(z)$  corresponding to each zone.

Due to the varying postures of passengers, the position of the speech source location within each zone can be arbitrary. When a speaker adopts a relatively standard sitting posture, it becomes easier to accurately separate each speaker's voice into the correct zone, which is called "standard postures". However, when a speaker's sitting posture is non-standard, their speech may occur at the boundary between adjacent zones, as indicated by the area within the red dotted line in Figure 1(a), which is called "non-standard postures". This may cause the speech separation system to incorrectly localize speech to adjacent zones.

## 3. Method

### 3.1. System overview

The overall architecture of our proposed CabinSep is shown in Figure 1(b). First, the  $Z$ -channel audio mixture  $y$  is transformed by the Short-Time Fourier Transform (STFT) to obtain the T-F spectrum  $Y \in \mathbb{C}^{Z \times T \times F}$ , where  $Z$  represents the channel dimension,  $T$  represents the time dimension, and  $F$  represents the frequency dimension. The log power spectrum (LPS)  $L$  and interaural phase difference (IPD)  $I$  [21, 29] are then derived from  $Y$ , which contain abundant spectral and spatial in-

formation, respectively. Next, we adopt three encoders to embed  $Y$ ,  $L$ , and  $I$ , respectively. These embeddings are concatenated along the channel dimension and passed through a convolution layer, resulting in  $E \in \mathbb{R}^{C \times T \times F}$ . Multiple full-sub modules are then adopted to refine these T-F embeddings [30], and the output is restored to  $Z \times T \times F$  through deconvolution. Finally, the mask estimate module is adopted to simultaneously estimate speech mask  $M_S$  and noise mask  $M_N$ . During inference, as shown in Figure 1(g), a streaming MVDR [31] utilizes the  $M_S$  and  $M_N$  along with  $Y$  as input to obtain the separated clean speech  $X_i$  for each zone:

$$X_{i,t,f} = W_{i,t,f}^H Y_{t,f}, \quad (3)$$

$$W_{i,t,f} = \frac{\Psi_{i,t,f}^{-1} \Phi_{i,t,f} e}{\text{tr}(\Psi_{i,t,f}^{-1} \Phi_{i,t,f})} \quad (4)$$

where  $W_{i,t,f}$  is the MVDR coefficient vector for each zone, and  $e$  is a one-hot vector with 1 at the reference microphone position. The matrices  $\Psi_i$  and  $\Phi_i$  represent the spatial covariance matrices of the target speech and overlapping interference, derived from  $M_S$  and  $M_N$ . The detailed calculation method is described in [11, 31].

### 3.2. Spec Encoder, LPS Encoder and IPD Encoder

The proposed CabinSep consists of three encoders. Each encoder consists of two convolution layers, and the ReLU is used as the activation after each convolution layer. The real and imaginary parts of  $Y$  are stacked along the channel dimension to obtain  $Y_R \in \mathbb{R}^{2Z \times T \times F}$ , which is then sent into the Spec encoder. The LPS encoder takes the  $L$  as input, which is defined as  $L = \log(|Y_R|^2) \in \mathbb{R}^{Z \times T \times F}$ . Due to the audio aliasing caused by large microphone spacing in the back-row (zone 3 and zone 4), the IPD information between these two microphones becomes unusable. Therefore, we only utilize the IPD information  $I$  between the front-row (zone 1 and zone 2) microphones with smaller microphone spacing. Following the method in [29],  $I$  is calculated by:

$$I = [\cos(\theta_{1,2}), \sin(\theta_{1,2})] \in \mathbb{R}^{2 \times T \times F}, \quad (5)$$

where  $\theta_{1,2} = \angle Y_1 - \angle Y_2$  represent the phases difference between the signals received by the microphones with index 1 and 2 in zone 1 and zone 2, respectively.

### 3.3. Full-Sub Modules

Embedded features are then refined by  $N$  full-sub modules. The value of  $N$  can be selected according to the limitations of parameters and computational complexity. Each full-sub module comprises three submodules: full-band-LSTM, TAC, and

Sub-band-conformer. The full-band-LSTM is used to process the full-band information, and its structure is the same as that in [30].

The structure of the TAC module is shown in Figure 1(c). Referring to [10, 19], TAC models spatial features across different channels. Considering that the computational complexity of the TAC module is relatively large, we propose a time skip operation before TAC. Specifically, we randomly select either the 1st or 2nd frame as the starting time frame, then choose frames to be processed at intervals of 1. The time skip operation halves the number of time frames, and the obtained feature  $\mathbf{E}' \in \mathbb{R}^{C \times \frac{T}{2} \times F}$  serves as the input for the TAC. Since TAC emphasizes information across channels and is insensitive to time frames, incorporating the time skip operation reduces the computational complexity of TAC by half with little impact on performance. In TAC, two fully connected layers,  $\text{Linear}_A$  and  $\text{Linear}_B$ , first reduce the dimension of the channel dimension of the input feature to obtain  $\mathbf{E}'_A$  and  $\mathbf{E}'_B$ :

$$\mathbf{E}' \in \mathbb{R}^{C \times \frac{T}{2} \times F} \rightarrow \mathbf{E}'_A \in \mathbb{R}^{\frac{C}{d} \times \frac{T}{2} \times F}, \mathbf{E}'_B \in \mathbb{R}^{\frac{C}{d} \times \frac{T}{2} \times F}, \quad (6)$$

where  $d$  is the compression ratio of the channel dimension. Then, the ReLU-activated  $\mathbf{E}'_B$  is averaged along the channel dimension and repeatedly stacked. It is then concatenated with the ReLU-activated  $\mathbf{E}'_A$  along the channel dimension to obtain  $\mathbf{E}'_{\text{cat}}$ , after which a fully connected layer  $\text{Linear}_C$  is used to restore the dimension. After TAC, the processed time frames selected by the time skip operation are recombined with the unprocessed time frames.

The sub-band-conformer models the features in sub-bands. Its structure is shown in Figure 1(d), and it is generally consistent with that of [30], except that the LSTM is replaced with the conformer [32, 33, 34, 35]. Due to the small number of parameters and low computational complexity associated with the sub-band processing structure, replacing the LSTM with the conformer introduces only a few additional parameters and has minimal impact on computational complexity, while significantly enhancing the temporal modeling ability.

### 3.4. Masks Estimate and Model Training

The mask estimate module, shown in Figure 1(e), estimates both the speech and noise masks using multiple linear layers. During training, the estimated speech and noise masks are applied to the multi-channel audio mixtures to obtain the separated speech and noise, and the losses are calculated with the corresponding labels. The loss function consists of  $L_{\text{Fbank-MAE}}$  and  $L_{\text{SI-SNR}}$  [22] where  $L_{\text{Fbank-MAE}}$  is the mean absolute error (MAE) loss calculated between the fbank feature of outputs and labels [36]. The composition of the overall loss function is:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{Fbank-MAE}}(\mathbf{S}, \mathbf{S}_{\text{label}}) + \beta \mathcal{L}_{\text{SI-SNR}}(\mathbf{S}, \mathbf{S}_{\text{label}}) + \gamma \mathcal{L}_{\text{Fbank-MAE}}(\mathbf{N}, \mathbf{N}_{\text{label}}), \quad (7)$$

where  $\mathbf{S}$ ,  $\mathbf{N}$  represent the separated speech and noise, and  $\mathbf{S}_{\text{label}}$  and  $\mathbf{N}_{\text{label}}$  are corresponding labels.  $\alpha$ ,  $\beta$  and  $\gamma$  represent the weight of those loss functions. We set  $\alpha = 0.01$ ,  $\beta = 1$  and  $\gamma = 0.01$  to balance the magnitude.

We adopt a two-stage training strategy to address inaccurate positioning caused by “non-standard postures” and further improve the speech recognition accuracy of back-end ASR models. In the first stage, we use simulated IRs for data augmentation. In the second stage, the model is finetuned by the augmented data using a mixture of simulated and real-recorded IRs.

## 4. Experiment

### 4.1. Datasets

**Training set:** The training set includes clean speech, background noise, transient noises including claps, coughs, etc., and IRs used to simulate reverberation. The clean speech is from AISHELL-2 [37] and augmented to simulate scenarios with one to four passengers speaking simultaneously in a car. The background and transient noise are recorded in cars and last 10 hours. The IRs include both simulated and real-recorded IRs. The simulated IRs are generated using the image-source method<sup>2</sup>, based on the actual size of the car cabins. We set the number of zones within the cabin to  $Z = 4$  and simulate 25,000 IRs for each zone’s microphone, resulting in a total of 100,000 simulated IRs. For generating the real-recorded IRs, we evaluate three different activation signals: exponential sine sweep signal (ESS) [38], maximum length sequence (MLS) [39], and time-stretched pulses (TSP) [40]. We record 39 IRs in each zone or each activation signal, leading to 156 IRs across the four zones.

**Test set:** Audios used for testing are recorded in real-world in-car scenarios consisting of two parts: the speech recognition test set and the zone positioning test set. The speech recognition test set lasts 7.4 hours, of which 37.5% of the recordings are recorded during driving, with a lower signal-to-noise ratio (SNR). The zone positioning test set lasts 4.9 hours, with 57% of the recordings made during car motion, 25% of the recordings having the windows opened, and 75% of the recordings involving “non-standard postures”. The remaining 43% of the scenes are recorded when the car is stationary, with 67% involving “non-standard postures”.

### 4.2. Data Augmentation

We use IRs and clean speech to simulate four-channel reverberant speech data, mimicking the recordings received by four microphones inside the cabin. In the 1st-stage of training, we convolve the clean speech with simulated IRs to simulate reverberant speech, while in the 2nd stage, we propose a data augmentation method that combines real-recorded IRs with simulated IRs. The channel corresponding to the speaker’s sitting zone uses real-recorded IRs to simulate reverberant speech, while the other three channels use simulated IRs to simulate reverberant speech. We name this method “mixed real-recorded IRs” in the experiment. Additionally, we compare two other methods for data simulation using real-recorded IRs. One method combines simulated IRs with real-recorded IRs in a certain ratio, where 25% of the data has all channels augmented with real-recorded IRs, and the remaining 75% uses simulated IRs. We call this method “added real-recorded IRs”. The other method only uses real-recorded IRs for reverberant speech simulation, which is named “only real-recorded IRs”. After simulating the four-channel reverberant data, we add noise in an on-the-fly manner during training. The SNR range of background noise is [-20, 25] db, while the SNR range of transient noise is [-5, 5] db.

### 4.3. Experimental Setup

We design three models with different sizes, namely CabinSep-S, CabinSep-M, and CabinSep-L, with the number of full-sub modules  $N$  being 1, 2, and 3, respectively. For CabinSep-S and CabinSep-M, the channel compression ratio  $d$  for  $\text{Linear}_A$  and  $\text{Linear}_B$  in the TAC is set to 4, while in CabinSep-L,  $d$  is set to 2. In CabinSep-S, the conformer consists of 4 layers, whereas the other two models

<sup>2</sup><https://github.com/DavidDiazGuerra/gpuRIR/>

Table 1: Comparison of CER on real-recorded speech recognition test set of the 1st-stage model trained with simulated IRs.

#	System	Causal	Params M	GMACs	WeNet [26]			SenseVoice [27]		
					static CER%	motion CER%	average CER%	static CER%	motion CER%	average CER%
1	Unprocessed Mixture	-	-	-	62.34	42.74	55.42	39.37	20.82	32.82
2	FasNet-TAC	✓	2.77	10.35	26.50	44.57	32.88	13.72	23.69	17.24
3	DualSep-S	✓	0.83	1.07	18.02	33.32	23.42	11.48	18.7	14.03
4	DualSep-L	✓	1.12	1.52	17.63	30.76	22.26	11.52	16.97	13.44
5	CabinSep-S	✓	1.09	0.40	16.87	22.09	18.36	10.77	13.44	11.53
6	CabinSep-M	✓	2.24	0.62	16.46	<b>20.96</b>	17.74	10.64	<b>13.07</b>	<b>11.33</b>
7	CabinSep-L	✓	3.43	1.22	<b>15.74</b>	21.48	<b>17.38</b>	<b>10.53</b>	13.28	<b>11.31</b>
7-1	-MVDR	✓	3.43	1.18	26.6	42.57	31.16	14.84	23.47	17.30
7-2	+time skip	✓	3.43	0.83	16.39	21.29	17.79	10.70	13.45	11.48
7-3	-conformer	✓	3.42	1.19	16.22	22.24	17.94	10.68	13.50	11.49
7-4	-conformer-TAC	✓	3.02	0.43	18.08	25.42	20.18	10.85	13.98	11.74
7-5	-conformer-TAC-LPS-IPD	✓	3.02	0.38	19.00	26.48	21.13	10.98	14.84	12.08
7-6	-conformer-TAC-Noise Mask	✓	2.83	0.38	19.02	26.35	21.11	10.98	14.55	12.00
7-7	+chunk	✓	3.43	1.22	15.96	21.26	17.47	10.56	13.27	11.33

Table 2: CER results obtained by WeNet after finetuning the 1st-stage CabinSep-L using real-recorded IRs. NSPA represents the positioning accuracy rate in “non-standard posture”.

Type	1st-stage IRs		mixed real-recorded IRs		added real-recorded IRs		only real-recorded IRs	
	CER%↓	NSPA%↑	CER%↓	NSPA%↑	CER%↓	NSPA%↑	CER%↓	NSPA%↑
ESS	17.38	60.4	16.61	95.4	16.77	98.9	17.26	98.9
MLS	17.38	60.4	16.68	91.7	16.77	93.9	17.08	95.2
TSP	17.38	60.4	16.59	93.4	16.64	98.1	16.79	97.6

have 2 layers each. In all three models, the convolution layer’s output dimension  $H$  of the sub-band block is 16, the feedforward dimension in the conformer is  $\frac{H}{2} = 8$ , and the multi-head attention uses 4 heads. We set the output dimension  $C$  of the Conv2d before full-sub modules to 24.

We adopt two other speech separation approaches, FasNet-TAC [19] and DualSep [20], as baselines for comparison. FasNet-TAC is a classic open source<sup>3</sup> end-to-end time-domain speech separation network, while DualSep is an in-car speech separation approach that reaches SOTA performance. We train and test these two models using the same datasets. For DualSep, we train the two models with different model sizes proposed in [20], namely DualSep-S and DualSep-L, for comparison. Since DualSep-L incorporates a non-causal IVA operation, making the model overall non-causal and unsuitable for practical applications, we substitute the non-causal IVA in DualSep-L with a causal IVA.

For both CabinSep and the DualSep systems, the input audio is transformed into the T-F domain via STFT. We use a hamming window with an FFT size of 512, window length of 32 ms, and hop length of 16 ms. During the training and finetuning, the Adam optimizer is adopted with an initial learning rate of 0.0001, which is halved every 20,000 steps. The separated speech is sent to two different back-end ASR models, WeNet<sup>4</sup> and SenseVoice<sup>5</sup>, to evaluate the character error rate (CER), respectively.

#### 4.4. Results

**Stage 1: Training on simulated IRs.** Table 1 shows the CER results of the 1st-stage systems on the speech recognition test set, of which the upper part compares our proposed CabinSep systems with baselines, and the lower part shows the results of ablation experiments. We use two different ASR models, WeNet [26] and SenseVoice [27], for speech recognition, comparing the static CER when the car is stationary, the motion CER when the vehicle is in motion, and the average CER on the whole test set. From the upper part of Table 1, we can see that CabinSep-S has significantly lower computational complexity than baselines while consistently achieving a lower CER. Specifically, CabinSep-S, with a computational complexity of only 0.4 GMACs and 0.21 RTF on a single-core Qualcomm SA8295P in-car CPU, achieves average CER relative reductions of 17.5% and 14.2% through WeNet and SenseVoice, respectively, compared with DualSep-L. With the larger model size and computational complex-

ity, CabinSep-M and CabinSep-L can receive a consistent improvement in CER performance. The best-performing model, CabinSep-L, achieved an average CER of 17.38% with WeNet and 11.31% with SenseVoice. Therefore, we selected CabinSep-L for the ablation experiments to evaluate the effectiveness of each of our proposed modules.

**Ablation Results:** Analyzing the average CER results from WeNet in the lower part of Table 1, it is evident that each proposed module is practical. The most significant improvement is observed with the cascaded MVDR (7-1), reducing CER from 31.16% to 17.38% with only an additional 0.04 GMACs. Introducing a time skip operation (7-2) reduces computational complexity by 0.39 GMACs, with a minimal CER increase of 0.41%. Replacing the conformer in the sub-band-conformer with LSTM (7-3) slightly decreases the parameter count and computational complexity by 0.01M and 0.03 GMACs, respectively, but increases CER by 0.56%. Further omitting TAC modules (7-4) results in a 2.24% increase in CER. Experiments 7-5 and 7-6 demonstrate that further removing the LPS and IPD encoders and estimating only the speech mask increases CER by 0.95% and 0.93%, respectively. Finally, adding chunks (7-7), limiting the conformer to look back at a maximum of 2 seconds of audio, only increases the CER by 0.09%.

**Stage 2: Finetuning on real-recorded IRs.** In the 2nd-stage experiments, we compare three types of activation signals to generate real-recorded IRs. For “standard postures”, all methods achieve 100% positioning accuracy. As shown in Table 2, real-recorded IRs significantly improve the positioning accuracy in “non-standard postures” (NSPA) scenarios, raising it from 60.4% to over 90%, with the highest accuracy reaching 98.9%. We found that IRs generated using ESS and TSP activation signals perform similarly and better than those generated by MLS. The lowest CER is achieved using our proposed “mixed real-recorded IRs” data augmentation method. Due to the limited amount of real-recorded IRs, CabinSep-L trained exclusively with real-recorded IRs for data augmentation performs worse in CER metric than using a combination of simulated and real-recorded IRs.

## 5. Conclusions

This paper proposes an in-car speech separation approach with excellent generalization ability, enhancing the back-end ASR models’ performance. We validate the system using two different ASR models, WeNet and SenseVoice, without any joint training, demonstrating its plug-and-play capability and compatibility with various ASR systems. Trained with simulated data and tested on real-recorded test sets, the proposed CabinSep-S achieves a 17.5% relative reduction in CER compared to the previous SOTA approach, DualSep-L, while requiring 1.12 GMACs less computational complexity. Additionally, we compared various data augmentation strategies that combine simulated and real-recorded IRs and addressed the issue of speaker mispositioning in-car scenarios, particularly for speakers in “non-standard postures”.

<sup>3</sup><https://github.com/yluo42/TAC/tree/master>

<sup>4</sup><https://github.com/wenet-e2e/wenet/blob/main/examples/wenetspeech/s0/>

<sup>5</sup><https://huggingface.co/FunAudioLLM/SenseVoiceSmall>

## 6. References

- [1] M. Capallera, L. Angelini, and et al., “Human-vehicle interaction to support driver’s situation awareness in automated vehicles: A systematic review,” *IEEE Trans. Intell. Veh.*, vol. 8, no. 3, pp. 2551–2567, 2023.
- [2] P. K. Murali, M. Kaboli, and et al., “Intelligent in-vehicle interaction technologies,” *Adv. Intell. Syst.*, vol. 4, no. 2, 2022.
- [3] K. Müller, S. Doclo, and et al., “Model-based estimation of in-car-communication feedback applied to speech zone detection,” in *IWAENC*. IEEE, 2022, pp. 1–5.
- [4] F. Weng, P. Angkititrakul, and et al., “Conversational in-vehicle dialog systems: The past, present, and future,” *IEEE Signal Process. Mag.*, vol. 33, no. 6, pp. 49–60, 2016.
- [5] W. Li, C. Miyajima, and et al., “Adaptive nonlinear regression using multiple distributed microphones for in-car speech recognition,” *IEICE Trans. Fundam. Electron. Commun. Comput. Sci.*, vol. 88-A, no. 7, pp. 1716–1723, 2005.
- [6] J. Barker, R. Marxer, and et al., “The third ‘chime’ speech separation and recognition challenge: Dataset, task and baselines,” in *ASRU*. IEEE, 2015, pp. 504–511.
- [7] J. Yu and S. Zhang, “Audio-visual multi-channel integration and recognition of overlapped speech,” *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 2067–2082, 2021.
- [8] P. Guo, H. Wang, and et al., “The NPU-ASLP system for audio-visual speech recognition in MISP 2022 challenge,” in *ICASSP*. IEEE, 2023, pp. 1–2.
- [9] B. Mu, P. Guo, and et al., “Automatic channel selection and spatial feature integration for multi-channel speech recognition across various array topologies,” in *ICASSP*. IEEE, 2024, pp. 11 396–11 400.
- [10] Y. Bando, T. Nakamura, and et al., “Neural blind source separation and diarization for distant speech recognition,” *arXiv preprint arXiv:2406.08396*, 2024.
- [11] T. Yoshioka, H. Erdogan, and et al., “Recognizing overlapped speech in meetings: A multichannel separation approach using neural networks,” in *Interspeech*, B. Yegnanarayana, Ed. ISCA, 2018, pp. 3038–3042.
- [12] J. Heymann, L. Drude, and et al., “Neural network based spectral mask estimation for acoustic beamforming,” in *ICASSP*. IEEE, 2016, pp. 196–200.
- [13] H. Erdogan, J. R. Hershey, and et al., “Improved MVDR beamforming using single-channel mask prediction networks,” in *Interspeech*, N. Morgan, Ed. ISCA, 2016, pp. 1981–1985.
- [14] Y. Xu, C. Weng, and et al., “Joint training of complex ratio mask based beamformer and acoustic model for noise robust asr,” in *ICASSP*. IEEE, 2019, pp. 6745–6749.
- [15] S. Zhao and D. L. Jones, “A fast-converging adaptive frequency-domain MVDR beamformer for speech enhancement,” in *INTER-SPEECH*. ISCA, 2012, pp. 1930–1933.
- [16] Z. Zhang, Y. Xu, and et al., “Multi-channel multi-frame ADL-MVDR for target speech separation,” *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 3526–3540, 2021.
- [17] Z. Wang, P. Wang, and et al., “Multi-microphone complex spectral mapping for utterance-wise and continuous speech separation,” *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 2001–2014, 2021.
- [18] A. Li, W. Liu, and et al., “Embedding and beamforming: All-neural causal beamformer for multichannel speech enhancement,” in *ICASSP*. IEEE, 2022, pp. 6487–6491.
- [19] Y. Luo, Z. Chen, and et al., “End-to-end microphone permutation and number invariant multi-channel speech separation,” in *ICASSP*. IEEE, 2020, pp. 6394–6398.
- [20] Z. Wang, J. Sun, and et al., “Dualsep: A light-weight dual-encoder convolutional recurrent network for real-time in-car speech separation,” in *SLT*. IEEE, 2024, pp. 286–293.
- [21] Y. Xu, V. Kothapally, and et al., “Zoneformer: On-device neural beamformer for in-car multi-zone speech separation, enhancement and echo cancellation,” in *Interspeech*. ISCA, 2023, pp. 5117–5121.
- [22] J. L. Roux, S. Wisdom, and et al., “SDR - half-baked or well done?” in *ICASSP*. IEEE, 2019, pp. 626–630.
- [23] R. Han, X. Yan, and et al., “An audio-quality-based multi-strategy approach for target speaker extraction in the misp 2023 challenge,” in *ICASSP-Workshops*. IEEE, 2024, pp. 27–28.
- [24] M. Delcroix, K. Zmolíková, and et al., “Single channel target speaker extraction and recognition with speaker beam,” in *ICASSP*. IEEE, 2018, pp. 5554–5558.
- [25] N. J. Bryan, “Impulse response data augmentation and deep neural networks for blind room acoustic parameter estimation,” in *ICASSP*. IEEE, 2020, pp. 1–5.
- [26] Z. Yao, D. Wu, and et al., “Wenet: Production oriented streaming and non-streaming end-to-end speech recognition toolkit,” in *Interspeech*. ISCA, 2021, pp. 4054–4058.
- [27] K. An, Q. Chen, and et al., “Funaudiollm: Voice understanding and generation foundation models for natural interaction between humans and llms,” *arXiv preprint arXiv:2407.04051*, 2024.
- [28] H. Wang, P. Guoand, and et al., “ICMC-ASR: the ICASSP 2024 in-car multi-channel automatic speech recognition challenge,” in *ICASSP-Workshops*. IEEE, 2024, pp. 63–64.
- [29] Z. Wang, J. L. Roux, and et al., “Multi-channel deep clustering: Discriminative spectral and spatial embeddings for speaker-independent speech separation,” in *ICASSP*. IEEE, 2018, pp. 1–5.
- [30] Z.-Q. Wang, S. Cornell, and et al., “Fneural speech enhancement with very low algorithmic latency and complexity via integrated full-and sub-band modeling,” in *ICASSP*. IEEE, 2023, pp. 1–5.
- [31] T. Higuchi, K. Kinoshita, and et al., “Frame-by-frame closed-form update for mask-based adaptive mvdr beamforming,” in *ICASSP*, 2018, pp. 531–535.
- [32] A. Gulati, J. Qin, and et al., “Conformer: Convolution-augmented transformer for speech recognition,” in *Interspeech*. ISCA, 2020, pp. 5036–5040.
- [33] S. Chen, Y. Wu, and et al., “Continuous speech separation with conformer,” in *ICASSP*. IEEE, 2021, pp. 5749–5753.
- [34] Y. Fu, Y. Liu, and et al., “Uformer: A unet based dilated complex & real dual-path conformer network for simultaneous speech enhancement and dereverberation,” in *ICASSP*. IEEE, 2022, pp. 7417–7421.
- [35] R. Han, W. Xu, and et al., “Distil-dccrn: A small-footprint DC-CRN leveraging feature-based knowledge distillation in speech enhancement,” *IEEE Signal Process. Lett.*, vol. 31, pp. 2075–2079.
- [36] T. Saramaki and R. Bregović, “Multirate systems and filter banks,” 2002. [Online]. Available: <https://api.semanticscholar.org/CorpusID:9768753>
- [37] J. Du, X. Na, and et al., “Aishell-2: Transforming mandarin asr research into industrial scale,” *arXiv preprint arXiv:1808.10583*, 2018.
- [38] A. Farina, “Simultaneous measurement of impulse response and distortion with a swept-sine technique,” in *Audio engineering society convention 108*. Audio Engineering Society, 2000.
- [39] H. Alrutz, “A fast hadamard transform method for the evaluation of measurements using pseudorandom test signals,” *Proc. 11th ICA, Paris, July 1983*, vol. 6, pp. 235–238, 1983.
- [40] N. Aoshima, “Computer-generated pulse signal applied for sound measurement,” *The Journal of the Acoustical Society of America*, vol. 69, no. 5, pp. 1484–1488, 1981.