# MFCCA:MULTI-FRAME CROSS-CHANNEL ATTENTION FOR MULTI-SPEAKER ASR IN MULTI-PARTY MEETING SCENARIO

*Fan Yu[1], Shiliang Zhang, Pengcheng Guo[1], Yuhao Liang[1], Zhihao Du, Yuxiao Lin[2], Lei Xie[1*]*

[1]Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science, Northwestern Polytechnical University, Xi'an, China
[2]College of Computer Science and Technology, Zhejiang University, Hangzhou, China

## ABSTRACT

Recently cross-channel attention, which better leverages multi-channel signals from microphone array, has shown promising results in the multi-party meeting scenario. Cross-channel attention focuses on either learning global correlations between sequences of different channels or exploiting fine-grained channel-wise information effectively at each time step. Considering the delay of microphone array receiving sound, we propose a multi-frame cross-channel attention, which models cross-channel information between adjacent frames to exploit the complementarity of both frame-wise and channel-wise knowledge. Besides, we also propose a multi-layer convolutional mechanism to fuse the multi-channel output and a channel masking strategy to combat the channel number mismatch problem between training and inference. Experiments on the AliMeeting, a real-world corpus, reveal that our proposed model outperforms single-channel model by 31.7% and 37.0% CER reduction on Eval and Test sets. Moreover, with comparable model parameters and training data, our proposed model achieves a new SOTA performance on the AliMeeting corpus, as compared with the top ranking systems in the ICASSP2022 M2MeT challenge, a recently held multi-channel multi-speaker ASR challenge.

*Index Terms*— Multi-speaker ASR, multi-channel, cross-channel attention, AliMeeting, M2MeT

## 1. INTRODUCTION

Multi-speaker automatic speech recognition (ASR) aims to transcribe speech that contains multiple speakers, and hopefully overlapped speech can be correctly transcribed. It is an essential task of rich transcription in multi-party meetings [1, 2, 3]. In recent years, with the advances of deep learning, many end-to-end neural multi-speaker ASR approaches have been proposed [4, 5, 6] and promising results have been achieved on synthetic multi-speaker datasets, e.g., LibriCSS [7]. However, transcribing real-world meetings is far more challenging with entangled difficulties such as overlapping speech, conversational speaking style, unknown number of speakers, far-field speech signals with noise and reverberation. Recently, two challenges – Multi-channel Multi-party Meeting Transcription (M2MeT) [8, 9] and Multimodal Information based Speech Processing (MISP) [10] – have made available valuable real-world multi-talker speech datasets to benchmark multi-speaker ASR towards real conditions and applications.

In the real-world applications, microphone array is usually adopted for far-field speech recording scenarios, including those in M2MET and MISP, where beamforming is a common algorithm to leverage spatial information for multi-channel speech enhancement. With the help of deep neural networks, time-frequency mask-based beamforming [11, 12, 13, 14, 13, 15] has shown superior performance in various multi-speaker benchmarks, such as AMI [16], CHiME [17, 18] and M2MeT [8, 9]. The mask estimation network needs to be trained with signal-level criteria on the simulated data where the reference speech is required. Simulated data has a clear gap with real-world data, and optimizing the signal-level criteria may not necessarily lead to lowered word error rate (WER) as well. Aiming to alleviate such mismatch, joint optimization of multi-channel front-end and ASR has been proposed [19, 20, 21, 22, 23, 24]. Under the joint learning framework, the whole system can be optimized with an ultimate ASR loss function by adopting real-world data without reference-cleaned signals.

The *attention* mechanism has been recently introduced to neural beamforming [24, 25], which performs recursive *nonlinear* beamforming on the data represented in a latent space. Specifically, *cross-channel attention* has been proposed to directly leverage multi-channel signals in a neural speech recognition system [26, 27]. Impressively, such an approach can bypass the complicated front-end formalization and integrate beamforming and acoustic modeling into an end-to-end neural solution. This cross-channel attention approach takes the frame-wise multi-channel signal as input and learns global correlations between sequences of different channels, which can be easily depicted as mapping each channel representation (query) with a set of channel-average representation

---

(key-value) pairs to an output [26, 27], namely **f**rame-**l**evel **c**ross-**c**hannel **a**ttention (**FLCCA**). Meanwhile, **c**hannel-**l**evel **c**ross-**c**hannel (**CLCCA**) attention has recently achieved remarkable performance on speech separation [28, 29] and speaker diarization [30, 31] tasks, even leading a system to win the first place in the speaker diarization track in M2MeT challenge [31]. Compared with FLCCA, CLCCA is computed along the channel dimension, the representations of each channel are combined with those of the other channels for each time step [28], which functions similarly as beamforming.

From our point of view, FLCCA and CLCCA can be complementary in capturing temporal and spatial information. Frame-level is less capable of extracting fine-grained channel-wise patterns since averaging the channel representations directly may deteriorate the individual channel information. Channel-level cross-channel attention, on the other hand, only focuses on leveraging spatial diversities and capturing inter-channel correlations on each time step, without considering the context relationship between different channels. Thus, in this paper, we exploit the complementarity between frame-level and channel-level cross-channel attention and propose a **m**ulti-**f**rame **c**ross-**c**hannel **a**ttention (**MFCCA**) by modeling both channel-wise and frame-wise information simultaneously. Direction of arrival (DOA) estimation [32] has been widely used for speech enhancement, which utilizes the delay of microphone array receiving the signal to estimate the sound source direction based on the phase difference. Inspired by the intuitive idea behind DOA, our proposed method will pay more attention to channel context between adjacent frames to model both frame-wise and channel-wise dependencies.

We build our MFCCA based multi-channel ASR within an attention based encoder-decoder (AED) structure [33]. Moreover, the multi-channel outputs from the encoder are aggregated by multi-layer convolution to reduce channel dimensions gradually. Although the cross-channel attention is independent of the number and geometry of microphones, it has the well-known performance degradation issue when number of microphones is reduced [30, 28]. In order to combat this issue, we propose a channel masking strategy. By randomly masking several channels from the original multi-channel input during training, our MFCCA approach becomes more stable and robust to the arbitrary number of channels.

To the best of our knowledge, we are the first to leverage cross-channel attention on a real meeting corpus – AliMeeting – to examine its ability in multi-speaker ASR in meeting scenarios. Experiments on the AliMeeting corpus show that our proposed multi-channel multi-speaker ASR model outperforms the single-channel multi-speaker ASR model by 31.7% and 37.0% relative CER reduction on Eval and Test sets, respectively. Moreover, with comparable model parameters and amount of the training data, our proposed model

achieves 16.1% and 17.5% CER on Eval and Test sets, which surpasses the best system in the M2MeT challenge, resulting in a new SOTA performance on the AliMeeting corpus.

## 2. FROM SINGLE-CHANNEL TO CROSS-CHANNEL ATTENTION

In this section, we first review the multi-headed self-attention commonly used in signal channel cases and then introduce the frame-level and channel-level cross-channel attentions, respectively. A single channel feature input is defined as $\mathbf{X}$, while a $C$-channel input is formulated as $\bar{\mathbf{X}} = [\mathbf{X}_0, \cdots, \mathbf{X}_{C-1}]$.

### 2.1. Single-channel attention

Single-channel attention, which is a standard self-attention structure, adopts the multi-headed scaled dot-product to learn the contextual information within a single channel of speech signal, as shown in Fig. 1a. The output of a single-channel attention for the $i$-th head is calculated as

$$
\begin{aligned}
\mathbf{Q}_i^{sc} &= \mathbf{X}\mathbf{W}_i^{sc,q} + (\mathbf{b}_i^{sc,q})^T \in \mathbb{R}^{T \times D}, \\
\mathbf{K}_i^{sc} &= \mathbf{X}\mathbf{W}_i^{sc,k} + (\mathbf{b}_i^{sc,k})^T \in \mathbb{R}^{T \times D}, \\
\mathbf{V}_i^{sc} &= \mathbf{X}\mathbf{W}_i^{sc,v} + (\mathbf{b}_i^{sc,v})^T \in \mathbb{R}^{T \times D}, \\
\mathbf{H}_i^{sc} &= \text{Softmax}\left(\frac{\mathbf{Q}_i^{sc}(\mathbf{K}_i^{sc})^T}{\sqrt{D}}\right)\mathbf{V}_i^{sc} \in \mathbb{R}^{T \times D},
\end{aligned}
\tag{1}
$$

where Softmax($\cdot$) is the column-wise softmax function, $\mathbf{W}_i^{sc,*}$ and $\mathbf{b}_i^{sc,*}$ are learnable weight and bias parameters for the $i$-th head respectively.

### 2.2. Frame-level cross-channel attention

Frame-level cross-channel attention [26, 27] learns not only the contextual information between time frames but also spatial information across channels, as shown in Fig. 1b. The $i$-th head of FLCCA is calculated as

$$
\begin{aligned}
\mathbf{Q}_i^{fl} &= \bar{\mathbf{X}}\mathbf{W}_i^{fl,q} + (\mathbf{b}_i^{fl,q})^T \in \mathbb{R}^{C \times T \times D}, \\
\mathbf{K}_i^{fl} &= \bar{\mathbf{X}}'\mathbf{W}_i^{fl,k} + (\mathbf{b}_i^{fl,k})^T \in \mathbb{R}^{C \times T \times D}, \\
\mathbf{V}_i^{fl} &= \bar{\mathbf{X}}'\mathbf{W}_i^{fl,v} + (\mathbf{b}_i^{fl,v})^T \in \mathbb{R}^{C \times T \times D}, \\
\mathbf{H}_i^{fl} &= \text{softmax}\left(\frac{\mathbf{Q}_i^{fl}(\mathbf{K}_i^{fl})^T}{\sqrt{D}}\right)\mathbf{V}_i^{fl} \in \mathbb{R}^{C \times T \times D},
\end{aligned}
\tag{2}
$$

$\bar{\mathbf{X}}' = [\bar{\mathbf{X}}_0', \cdots, \bar{\mathbf{X}}_{C-1}']$. $\bar{\mathbf{X}}_c'$ is the average of all channels except for the $c^{th}$ channel, which is calculated by $\bar{\mathbf{X}}_c' = (\sum_{n, n \neq c} \bar{\mathbf{X}}_n)/(C-1) \in \mathbb{R}^{T \times D}$. $\mathbf{W}^{fl,*}$ and $\mathbf{b}^{fl,*}$ are learnable weight and bias parameters, respectively.
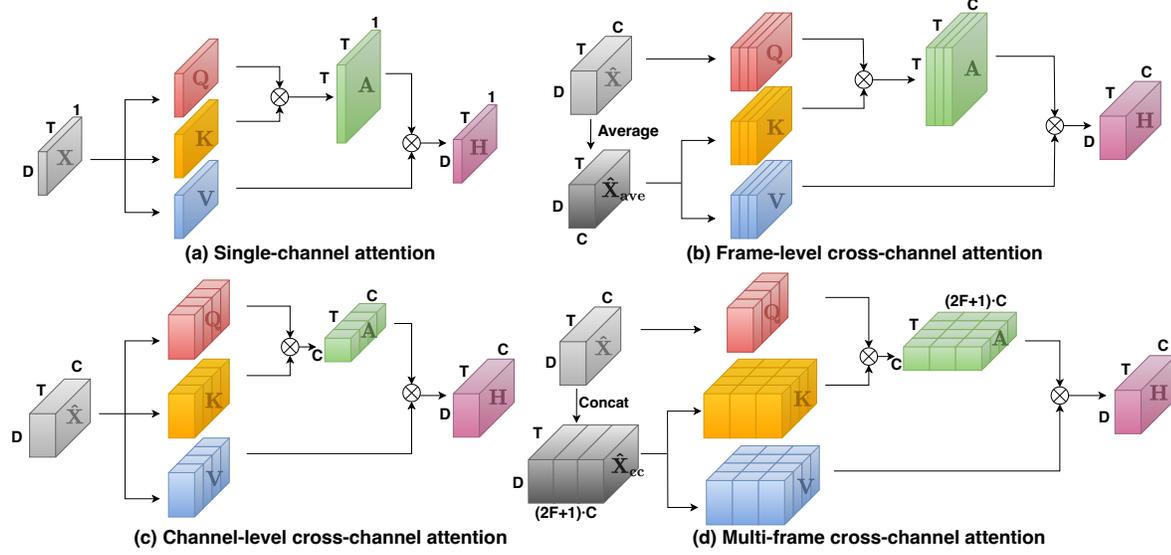
**Fig. 1.** Illustration of different attention blocks: (a) Single-channel attention. (b) Frame-level cross-channel attention (FLCCA). (c) Channel-level cross-channel attention (CLCCA). (d) Multi-frame cross-channel attention (MFCCA, proposed).

## 2.3. Channel-level cross-channel attention

Channel-level cross-channel attention focuses on leveraging spatial diversities and capturing inter-channel correlations on each time step, as shown in Fig. 1c. The $i$-th head of CLCCA can be formulated as

$$
\begin{aligned}
\mathbf{Q}_i^{cl} &= \bar{\mathbf{X}}\mathbf{W}_i^{cl,q} + (\mathbf{b}_i^{cl,q})^T \in \mathbb{R}^{T\times C\times D}, \\
\mathbf{K}_i^{cl} &= \bar{\mathbf{X}}\mathbf{W}_i^{cl,k} + (\mathbf{b}_i^{cl,k})^T \in \mathbb{R}^{T\times C\times D}, \\
\mathbf{V}_i^{cl} &= \bar{\mathbf{X}}\mathbf{W}_i^{cl,v} + (\mathbf{b}_i^{cl,v})^T \in \mathbb{R}^{T\times C\times D}, \\
\mathbf{H}_i^{cl} &= \text{softmax}\left(\frac{\mathbf{Q}_i^{cl}(\mathbf{K}_i^{cl})^T}{\sqrt{D}}\right)\mathbf{V}_i^{cl} \in \mathbb{R}^{T\times C\times D},
\end{aligned}
\tag{3}
$$

Again, $\mathbf{W}^{cl,*}$ and $\mathbf{b}^{cl,*}$ are learnable weight and bias parameters, respectively.

# 3. PROPOSED METHOD

## 3.1. Multi-frame cross-channel attention

Based on the discussion of FLCCA and CLCCA, multi-frame cross-channel attention is proposed to exploit the complementarity between frame-level and channel-level information, as shown in Fig. 1d. The $i$-th head of MFCCA is calculated as

$$
\begin{aligned}
\mathbf{Q}_i^{mf} &= \bar{\mathbf{X}}\mathbf{W}_i^{mf,q} + (\mathbf{b}_i^{mf,q})^T \in \mathbb{R}^{T\times C\times D}, \\
\mathbf{K}_i^{mf} &= \bar{\mathbf{X}}_{cc}\mathbf{W}_i^{mf,k} + (\mathbf{b}_i^{mf,k})^T \in \mathbb{R}^{T\times(2F+1)\cdot C\times D}, \\
\mathbf{V}_i^{mf} &= \bar{\mathbf{X}}_{cc}\mathbf{W}_i^{mf,v} + (\mathbf{b}_i^{mf,v})^T \in \mathbb{R}^{T\times(2F+1)\cdot C\times D}, \\
\mathbf{H}_i^{mf} &= \text{softmax}\left(\frac{\mathbf{Q}_i^{mf}(\mathbf{K}_i^{mf})^T}{\sqrt{D}}\right)\mathbf{V}_i^{mf} \in \mathbb{R}^{T\times C\times D},
\end{aligned}
\tag{4}
$$

where $\mathbf{W}^{mf,*}$ and $\mathbf{b}^{mf,*}$ are learnable weight and bias parameters, $\bar{\mathbf{X}}_{\mathbf{cc}} = [\bar{\mathbf{X}}_{\mathbf{cc}}^{\mathbf{0}}, \cdots, \bar{\mathbf{X}}_{\mathbf{cc}}^{\mathbf{t}}, \cdots, \bar{\mathbf{X}}_{\mathbf{cc}}^{\mathbf{T}}]$. $\bar{\mathbf{X}}_{cc}^t$ is the concatenation of the context frames, which is calculated by $\bar{\mathbf{X}}_{cc}^t =$

$[\bar{\mathbf{X}}^{t-F}, ..., \bar{\mathbf{X}}^t, ..., \bar{\mathbf{X}}^{t+F}] \in \mathbb{R}^{(2F+1)\cdot C\times D}$. $F$ is the number of the past and future frames to be concatenated at each time step, which is a trade-off between performance and computation cost. Inspired by the DOA calculation which utilizes the delay of the microphone array to estimate the source direction for speech enhancement, our proposed MFCCA focuses on channel context of adjacent frames to improve the ability of modeling the frame-level and channel-level contextual information together.

## 3.2. Conformer block

Our encoder layer also adopts the Conformer block [34, 35], which includes a multi-headed self-attention (MHSA) module, a convolution (CONV) module, and a pair of feed-forward (FFN) module in the Macaron-Net style. Conformer models both local and global dependencies of the audio sequence in a parameter-efficient way, which makes full use of the long-range global modeling ability of the MHSA module and the fine-grained local feature extraction ability of the CONV module. Note that the CONV and FFN module directly follow the multi-frame cross-channel attention will determinate the model performance, which will bring about 1% absolute CER reduction according to our experiment. Since the CONV module and FFN module both models at the frame-level, learning of channel dependence by multi-frame cross-channel attention will be affected. Thus, we adopt the model structure in Fig. 2.

## 3.3. Convolution fusion

To integrate the multi-channel outputs, previous studies [27, 31] mostly averaged or concatenated channel features along the time axis. In order to mitigate the corruption of channel-specific information caused by reducing the channel dimensions directly, we use a multi-layer convolution module to re-
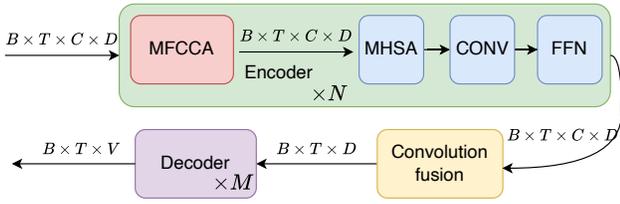
**Fig. 2**. An overview of the proposed multi-channel transformer network.

duce the channel dimensions gradually. As show in Fig. 3, the multi-layer convolution module consists of five 2-D convolution layers, which only increases negligible parameters. The number of input channels in the multi-layer convolution module is fixed. Therefore, if the channel number of the input is less than the pre-configured value, we need to expand the channel by simple repeating.
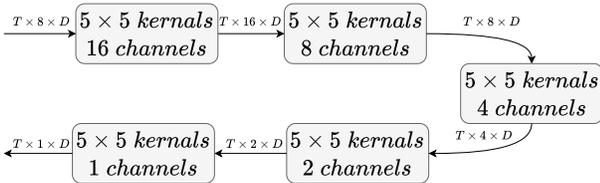


**Fig. 3**. The architecture of multi-layer convolution module.

### 3.4. Channel masking

Cross-channel attention is independent of the number of microphones and microphone geometry in its nature. But in practice, the performance of channel-level cross-channel attention is easily affected by the number of the channels [30, 28], especially when the channel numbers involved in the decoding and training period are different. Channel dropout [30] was proposed to prevent the models from being overly dependent on spatial information, in which multi-channel inputs are randomly dropped to be a single channel. However, channel dropout mainly improves the speech recognition performance of multi-channel model on a single-channel test set, which does not completely solve the problem of channel number mismatch. In order to improve the robustness of the model for different channel numbers, we introduce a channel masking strategy, which masks channels randomly for the multi-channel input. Specifically, a uniform probability $p \in (0, 1)$ is used to decide whether the multi-channel input will be masked. When choosing to mask, we randomly select $m \in (1, C)$ channels to be masked where $C$ is the total number of channels and $m$ is determined with equal probability $\frac{1}{C}$. Based on the channel masking strategy, our multi-channel ASR model can easily generalize to variant channel numbers as well as different microphone array geometries, leading to a more practical solution.

### 3.5. Training strategy

Considering the problem of overlapping speech and unknown number of speakers in real-world meeting scenarios, we adopt the Serialized Output Training (SOT) [6] to enable the multi-speaker recognition ability. The SOT scheme gets rid of the limitation on the number of speakers and models the dependencies among outputs of different speakers in an effective and simple way. In the training period, transcriptions of different speakers are serialized into a single word sequence with a special token ⟨sc⟩ inserted. The order of the transcriptions is determined by their start time. The experiments have shown that the SOT scheme achieves a better CER than the permutation invariant training (PIT) scheme, which needs to calculate all the permutations [6].

## 4. EXPERIMENTS

### 4.1. Dataset

We use AliMeeting[1] corpus [8, 9], a challenging Mandarin meeting dataset with multi-talker conversations, to evaluate our multi-channel multi-speaker ASR model. The AliMeeting corpus contains 104.75 hours data for training (Train), 4 hours for evaluation (Eval), and 10 hours for test (Test). Each set contains several meeting sessions and each session consists of a 15 to 30 minutes discussion by 2 to 4 participants. The AliMeeting corpus contains the 8-channel audios recorded from an annular microphone array (*Ali-far*), as well as the near-field audio (*Ali-near*) from the participant's headset microphone. *Ali-far-bf* is produced by applying CDDMA Beamformer [36, 37]. Meanwhile, similar to the M2MeT challenge submissions [38], we also use the training set of the *Aishell4*[2] [39] and 600 hours simulated training data named *Ali-simu* from *Ali-near*, which covers 2-4 speakers in one utterance with 15-40% overlapping ratio.

### 4.2. Baselines

We compare our MFCCA based multi-channel multi-speaker ASR model with four baselines: (1) *Single channel model*: as the single channel baseline. Specifically, we use the first channel of *Train-Ali-far* for training and testing. (2) *Beamformer*: the CDDMA Beamformer [36, 37] has shown promising results in speech enhancement and it uses all the channels for beamforming, which generates enhanced single channel data (*Ali-far-bf*) for the ASR model. (3) *Random selection*: a dynamic strategy is adopted to randomly select a channel of *Train-Ali-far* as the input to the ASR model during training. Note that the first channel is selected as the input for testing. (4) *Complex convolution*: the multi-channel real and imaginary parts of Short-Time Fourier Transform (STFT) results are extracted for complex convolution [40]. The convolution structure is similar to that in Fig. 3.

### 4.3. Experimental setup

In all experiments, we use the 80-dimensional Mel-filterbank feature extracted with a 25 ms frame length and a 10 ms window shift. The ESPnet [41] toolkit is used to build all our

---

**Table 1**. Results for various multi-channel approaches on Eval and Test sets (%).

| Model | Eval | | | | | Test | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 1-ch | 2-ch | 4-ch | 6-ch | 8-ch | 1-ch | 2-ch | 4-ch | 6-ch | 8-ch |
| Single channel [8, 9] | 32.3 | 32.3 | 32.3 | 32.3 | 32.3 | 33.8 | 33.8 | 33.8 | 33.8 | 33.8 |
| Beamformer [8, 9] | - | - | - | - | 30.7 | - | - | - | - | 31.8 |
| Random select | 30.2 | 30.2 | 30.2 | 30.2 | 30.2 | 31.2 | 31.2 | 31.2 | 31.2 | 31.2 |
| Complex convolution | 56.3 | 35.8 | 33.0 | 32.4 | 30.1 | 55.6 | 38.4 | 34.7 | 32.4 | 31.0 |
| Frame-level cross-channel [26, 27][†] | 60.5 | 50.4 | 25.9 | 22.6 | 22.5 | 63.8 | 51.8 | 27.5 | 24.6 | 24.6 |
| Channel-level cross-channel [30, 31][†] | 38.4 | 27.7 | 21.5 | 20.8 | 20.6 | 39.3 | 29.3 | 23.2 | 22.7 | 22.4 |
| Frame-level co-attention [30][†] | 38.1 | 26.3 | 23.2 | 22.7 | 22.5 | 39.1 | 27.9 | 24.4 | 24.2 | 24.0 |
| Multi-frame cross-channel | 38.0 | 27.3 | 21.2 | 20.6 | 20.2 | 39.0 | 28.8 | 22.9 | 22.3 | 22.0 |
| + Convolution fusion | 37.8 | 26.9 | 20.8 | 20.1 | 19.9 | 38.8 | 28.5 | 22.6 | 22.1 | 21.8 |
| + Mask channel ($p$=10%) | 36.1 | 25.8 | 20.3 | 19.7 | 19.6 | 37.2 | 27.6 | 22.2 | 21.8 | 21.5 |
| + Mask channel ($p$=15%) | 35.5 | 25.5 | 20.0 | 19.5 | 19.4 | 36.8 | 27.3 | 22.2 | 21.6 | 21.4 |
| + Mask channel ($p$=20%) | **35.1** | 25.4 | **20.0** | **19.5** | **19.4** | **36.3** | **26.9** | **22.0** | **21.5** | **21.3** |
| + Mask channel ($p$=25%) | 35.2 | **25.3** | 20.2 | 19.6 | 19.5 | 36.6 | 27.7 | 22.1 | 21.6 | 21.4 |

†: This models is re-implemented by ourselves with the same parameter structure as our model for fair comparison.

ASR systems. We follow the standard configuration of ES-Pnet to train the baseline models, which contain a 12-layer encoder and 6-layer decoder. The dimension of MHSA and FFN layers are set to 256 and 2048, respectively. For the cross-channel based models, we use an 11-layer encoder and a 6-layer decoder with the 4-head MHSA instead, in order to achieve a similar parameter size to the baseline models. All the ASR models are trained for 100 epochs and a warmup of the learning rate is used for the first 25,000 iterations. We use 4950 commonly used Mandarin characters as the modeling units. Results of all the experiments are measured by Character Error Rate (CER).

### 4.4. Comparison of different multi-channel models

As shown in Table 1, our proposed MFCCA model outperforms the four baselines, especially for the single channel model, leading to 31.7% (32.3%→19.4%) and 37.0% (33.8%→21.3%) relative CER reduction on 8-ch Eval and Test sets, respectively. Compared with other multi-channel attention models, our MFCCA model shows superior performance, achieving the lowest CER of 20.2% and 22.0% on 8-ch Eval and Test sets. When incorporating with the multi-layer convolution fusion to integrate multiple channels, we can obtain further improvement, decreasing the CER from 20.2%/22.0% to 19.9%/21.8% on 8-ch Eval and Test sets, respectively.

Cross-channel attention models perform well when the channel number of test set is large, but degrade significantly when the number of channels is reduced, e.g., single channel and 2-ch Test sets. Channel masking can improve the robustness of the model with different channel setups. According to the results, our model obtains the best results on most test sets when channel masking probability is set to 20%, achieving 7.1% (37.8%→35.1%) and 6.4% (38.8%→36.3%) relative CER reduction on 1-ch Eval and Test sets. Meanwhile,

channel masking also improves the multi-channel test sets and achieves CERs of 19.4% and 21.3% on 8-ch Eval and Test sets, which even has surpassed most of the submissions in the M2MeT challenge [8, 9].

### 4.5. Impact of the context frame number

As shown in Table 2, $F$ is the number of frames that looks back to the past and looks ahead to the future at each time step. When increasing $F$ from 0 to 2, we observe that the CER is improved from 20.6% to 20.0% on Eval set and 22.4% to 22.0% on Test set. When further increasing the $F$ from 2 to 4, the gain is marginal on Eval and Test sets, which only brings 0.1% absolute CER reduction on Eval set. The reason might be that the channel information of adjacent frames is more important in cross-channel attention, which denotes the importance of the delay time between microphones. Based on this conclusion, the frame number for looking back and ahead is set to 2 in the remaining experiments.

**Table 2**. Results of MFCCA model with different context frame number on Eval and Test sets (%).

| F | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| Eval | 20.6 | 20.4 | 20.2 | 20.2 | **20.1** |
| Test | 22.4 | 22.1 | **22.0** | 22.0 | 22.0 |

### 4.6. Visualization of MFCCA scores

To analyze the behavior of our proposed model, Fig. 4 visualizes the attention scores of our MFCCA module and the detailed recording process of the microphone array. As shown in Fig. 4(e), the different microphone-speaker distances may result in time delays during the recording. For example, the 7-th channel of speaker-1 shows a slight time delay compared with the 4-th channel, since the 4-th microphone is much closer to the speaker. Fig. 4 (a-d) are heatmaps of the averaged attention scores computed by our MFCCA module for different

speakers. As described in 3.1, for a specific time $t$, its input feature will be appended with two past and future contexts, and the MFCCA module tries to exploit cross-channel dependencies between adjacent frames. Combining Fig. 4(a) and Fig. 4(e), we can find that our model indeed captures the microphone delay information like beamforming, as the model attends more on the 4-th/5-th channels at time $t-2$ and 7-th channel at time $t$. Note that the attention scores are from the first encoder layer, in which each channel has not yet integrated the other channel information.



(a) Spkr 1 attention score  (b) Spkr 2 attention score  (c) Spkr 3 attention score  (d) Spkr 4 attention score

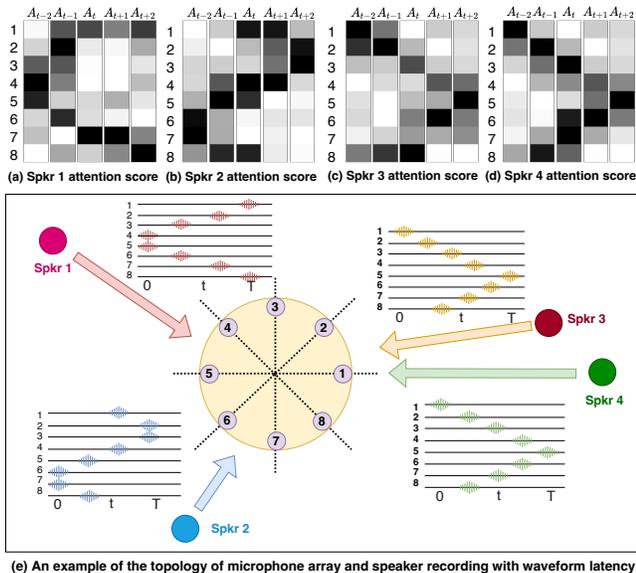(e) An example of the topology of microphone array and speaker recording with waveform latency

**Fig. 4**. Illustration of: (a-d) the attention scores of different speakers. (e) An example of the topology of microphone array and the recorded 8-ch waveforms.

### 4.7. Impact of the different training data scale

As shown in Table 3, we compare the results of our proposed model trained with different data scales on Eval and Test sets. In order to strengthen the acoustic modeling ability of the model, we include the *Train-Ali-near* and *Aishell4* sets into our training, which yields 10.8% (19.4%→17.3%) and 13.6% (21.3%→18.4%) relative CER reductions on Eval and Test sets, respectively. We also simulate 600 hours 8-channel meeting data based on the *Train-Ali-near* to have a fair comparison with the M2MeT challenge submissions. By using the same simulated data augmentation strategy, our model obtains further improvement, achieving 16.5% and 18.0% CERs on Eval and Test sets. Meanwhile, we also integrate a neural network language model (NNLM) into our proposed model to improve the language generalization ability, which brings 2.4% (16.5%→16.1%) and 2.7% (18.0%→17.5%) relative CER reductions on Eval and Test sets. The NNLM is trained on the transcriptions of training data, using extra text data is prohibited according to the M2MeT challenge rule.

Compared with the submission system of the $2^{nd}$ ranking team in M2MeT, which adopted the front-end and back-end joint modeling scheme [9, 38], our proposed MFCCA model brings 16.1% (19.2%→16.1%) and 15.9% (20.8%→17.5%) relative CER reductions on Eval and Test sets, while the parameters and training data are at a comparable scale. Furthermore, our model even outperforms the large model of the $1^{st}$ ranking team's submission system [9, 42] trained on a large data scale by data augmentation and simulation, leading to 8.0% (17.5%→16.1%) and 6.9% (18.8%→17.5%) relative CER reductions on Eval and Test sets, respectively.

**Table 3**. Results of MFCCA model with the different training data scales on Eval and Test sets (%).

| Model | Para(M) | Data(hrs) | Eval | Test |
|---|---|---|---|---|
| $1^{st}$ranking w/ model fusion[42] | 114 | 14,000 | 17.5 | 18.8 |
| $1^{st}$ranking [42] | 114 | 10,000 | 19.1 | 20.1 |
| $2^{nd}$ranking [38] | 48 | 917 | 19.2 | 20.8 |
| MFCCA (*Train-Ali-far*) | 45 | 105 | 19.4 | 21.3 |
| + *Train-Ali-near*, *Aishell4* | 45 | 317 | 17.3 | 18.4 |
| + *Ali-simu* | 45 | 917 | 16.5 | 18.0 |
| + NNLM | 45 | 917 | **16.1** | **17.5** |

## 5. CONCLUSIONS

In this work, we propose a multi-frame cross-channel attention (MFCCA) module based on the multi-speaker SOT framework to capture both temporal and spatial information, which exploits the complementarity between frame-level and channel-level cross-channel attention. Considering the delay of microphone array receiving sound, our MFCCA approach models cross-channel information between adjacent frames. Besides, we also propose a multi-layer convolutional mechanism to fuse the multi-channel output efficiently. Finally, in order to combat the channel number mismatch problem between training and inference, we propose a channel masking strategy to improve the robustness of the model with different channel setups. Evaluated on the real meeting corpus AliMeeting, our proposed model outperforms single channel ASR model by 31.7% and 37.0% relative CER reductions on Eval and Test sets, respectively. Moreover, with the comparable model parameters and training data, our proposed model achieves a SOTA error rate compared with top ranking systems in the ICASSP2022 M2MeT challenge, the recently held multi-channel multi-speaker ASR challenge. In the future, we would like to integrate our proposed multi-channel multi-speaker model into speaker-attributed automatic speech recognition for real-world applications.

# 7. REFERENCES

[1] Jonathan G Fiscus, Nicolas Radde, John S Garofolo, Audrey Le, Jerome Ajot, and Christophe Laprun, "The rich transcription 2005 spring meeting recognition evaluation," in *Proc. MLMI*. Springer, 2005, pp. 369–389.

[2] Jonathan G Fiscus, Jerome Ajot, Martial Michel, and John S Garofolo, "The rich transcription 2006 spring meeting recognition evaluation," in *Proc. MLMI*. Springer, 2006, pp. 309–322.

[3] Jonathan G Fiscus, Jerome Ajot, and John S Garofolo, "The rich transcription 2007 meeting recognition evaluation," in *Proc. MTPH*, pp. 373–389. Springer, 2007.

[4] Dong Yu, Xuankai Chang, and Yanmin Qian, "Recognizing multi-talker speech with permutation invariant training," in *Proc. INTERSPEECH*. ISCA, 2017, pp. 2456–2460.

[5] Zhehuai Chen, Jasha Droppo, Jinyu Li, and Wayne Xiong, "Progressive joint modeling in unsupervised single-channel overlapped speech recognition," *Proc. TASLP*, vol. 26, no. 1, pp. 184–196, 2017.

[6] Naoyuki Kanda, Yashesh Gaur, Xiaofei Wang, Zhong Meng, and Takuya Yoshioka, "Serialized output training for end-to-end overlapped speech recognition," in *Proc. INTERSPEECH*. ISCA, 2020, pp. 2797–2801.

[7] Zhuo Chen, Takuya Yoshioka, Liang Lu, Tianyan Zhou, Zhong Meng, Yi Luo, Jian Wu, and Jinyu Li, "Continuous speech separation: dataset and analysis," in *Proc. ICASSP*. IEEE, 2020, pp. 7284–7288.

[8] Fan Yu, Shiliang Zhang, Yihui Fu, Lei Xie, Siqi Zheng, Zhihao Du, et al., "M2MeT: The ICASSP 2022 multi-channel multi-party meeting transcription challenge," in *Proc. ICASSP*. IEEE, 2022, pp. 6167–6171.

[9] Fan Yu, Shiliang Zhang, Pengcheng Guo, Yihui Fu, Zhihao Du, Siqi Zheng, Lei Xie, et al., "Summary on the ICASSP 2022 multi-channel multi-party meeting transcription grand challenge," in *Proc. ICASSP*. IEEE, 2022, pp. 9156–9160.

[10] Hang Chen, Hengshun Zhou, Jun Du, Chin-Hui Lee, Jingdong Chen, Shinji Watanabe, Sabato Marco Siniscalchi, Odette Scharenborg, Di-Yuan Liu, Bao-Cai Yin, Jia Pan, Jian-Qing Gao, and Cong Liu, "The first multimodal information based speech processing (MISP) challenge: Data, tasks, baselines and results," in *Proc. ICASSP*. IEEE, 2022, pp. 9266–9270.

[11] DeLiang Wang and Jitong Chen, "Supervised speech separation based on deep learning: An overview," in *Proc. TASLP*. 2018, vol. 26, pp. 1702–1726, IEEE.

[12] Keisuke Kinoshita, Marc Delcroix, Haeyong Kwon, Takuma Mori, and Tomohiro Nakatani, "Neural network-based spectrum estimation for online wpe dereverberation.," in *Proc. INTERSPEECH*. ISCA, 2017, pp. 384–388.

[13] Jahn Heymann, Lukas Drude, and Reinhold Haeb-Umbach, "Neural network based spectral mask estimation for acoustic beamforming," in *Proc. ICASSP*. IEEE, 2016, pp. 196–200.

[14] Dong Yu, Morten Kolbæk, Zheng-Hua Tan, and Jesper Jensen, "Permutation invariant training of deep models for speaker-independent multi-talker speech separation," in *Proc. ICASSP*. IEEE, 2017, pp. 241–245.

[15] Hakan Erdogan, John R Hershey, Shinji Watanabe, Michael I Mandel, and Jonathan Le Roux, "Improved mvdr beamforming using single-channel mask prediction networks.," in *Proc. INTERSPEECH*. ISCA, 2016, pp. 1981–1985.

[16] Iain McCowan, Jean Carletta, Wessel Kraaij, Simone Ashby, S Bourban, M Flynn, et al., "The AMI meeting corpus," in *Proc. ICMT*. Citeseer, 2005, vol. 88, p. 100.

[17] Jon Barker, Shinji Watanabe, Emmanuel Vincent, and Jan Trmal, "The fifth 'CHiME' speech separation and recognition challenge: Dataset, task and baselines," in *Proc. INTERSPEECH*. ISCA, 2018, pp. 1561–1565.

[18] Shinji Watanabe, Michael Mandel, Jon Barker, et al., "CHiME-6 Challenge: Tackling Multispeaker Speech Recognition for Unsegmented Recordings," in *Proc. CHiME 2020*, 2020, pp. 1–7.

[19] Naoyuki Kanda, Rintaro Ikeshita, Shota Horiguchi, Yusuke Fujita, Kenji Nagamatsu, Xiaofei Wang, Vimal Manohar, Nelson Enrique Yalta Soplin, Matthew Maciejewski, Szu-Jui Chen, et al., "The hitachi/jhu chime-5 system: Advances in speech recognition for everyday home environments using multiple microphone arrays," in *Proc. CHiME-5*, 2018, pp. 6–10.

[20] Naoyuki Kanda, Yusuke Fujita, Shota Horiguchi, Rintaro Ikeshita, Kenji Nagamatsu, and Shinji Watanabe, "Acoustic modeling for distant multi-talker speech recognition with single-and multi-channel branches," in *Proc. ICASSP*. IEEE, 2019, pp. 6630–6634.

[21] Aswin Shanmugam Subramanian, Xiaofei Wang, Murali Karthick Baskar, Shinji Watanabe, Toru Taniguchi, Dung Tran, and Yuya Fujita, "Speech enhancement using end-to-end speech recognition objectives," in *Proc. WASPAA*. IEEE, 2019, pp. 234–238.

[22] Aswin Shanmugam Subramanian, Chao Weng, Meng Yu, Shi-Xiong Zhang, Yong Xu, Shinji Watanabe, and Dong Yu, "Far-field location guided target speech extraction using end-to-end speech recognition objectives," in *Proc. ICASSP*. IEEE, 2020, pp. 7299–7303.

[23] Wangyou Zhang, Aswin Shanmugam Subramanian, Xuankai Chang, Shinji Watanabe, and Yanmin Qian, "End-to-end far-field speech recognition with unified dereverberation and beamforming," in *Proc. INTERSPEECH*. ISCA, 2020, pp. 324–328.

[24] Xuankai Chang, Wangyou Zhang, Yanmin Qian, Jonathan Le Roux, and Shinji Watanabe, "End-to-end multi-speaker speech recognition with transformer," in *Proc. ICASSP*. IEEE, 2020, pp. 6134–6138.

[25] Bahareh Tolooshams, Ritwik Giri, Andrew H Song, Umut Isik, and Arvindh Krishnaswamy, "Channel-attention dense u-net for multichannel speech enhancement," in *Proc. ICASSP*. IEEE, 2020, pp. 836–840.

[26] Feng-Ju Chang, Martin Radfar, Athanasios Mouchtaris, Brian King, and Siegfried Kunzmann, "End-to-end multi-channel transformer for speech recognition," in *Proc. ICASSP*. IEEE, 2021, pp. 5884–5888.

[27] Feng-Ju Chang, Martin Radfar, Athanasios Mouchtaris, and Maurizio Omologo, "Multi-channel transformer transducer for speech recognition," 2021.

[28] Dongmei Wang, Zhuo Chen, and Takuya Yoshioka, "Neural speech separation using spatially distributed microphones," in *Proc. INTERSPEECH*. ISCA, 2020, pp. 339–343.

[29] Dongmei Wang, Takuya Yoshioka, Zhuo Chen, Xiaofei Wang, Tianyan Zhou, and Zhong Meng, "Continuous speech separation with ad hoc microphone arrays," in *Proc. EUSIPCO*. IEEE, 2021, pp. 1100–1104.

[30] Shota Horiguchi, Yuki Takashima, Paola Garcia, Shinji Watanabe, and Yohei Kawaguchi, "Multi-channel end-to-end neural diarization with distributed microphones," in *Proc. ICASSP*. IEEE, 2022, pp. 7332–7336.

[31] Weiqing Wang, Xiaoyi Qin, and Ming Li, "Cross-channel attention-based target speaker voice activity detection: Experimental results for the m2met challenge," in *Proc. ICASSP*. IEEE, 2022, pp. 9171–9175.

[32] Jahn Heymann, Lukas Drude, Reinhold Haeb-Umbach, Keisuke Kinoshita, and Tomohiro Nakatani, "Frame-online dnn-wpe dereverberation," in *proc. IWAENC*. IEEE, 2018, pp. 466–470.

[33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Proc. NeurIPS*, 2017, pp. 5998–6008.

[34] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, et al., "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. INTERSPEECH*. ISCA, 2020, pp. 5036–5040.

[35] Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, et al., "Recent developments on ESPnet toolkit boosted by Conformer," in *Proc. ICASSP*. IEEE, 2021, pp. 5874–5878.

[36] Weilong Huang and Jinwei Feng, "Differential beamforming for uniform circular array with directional microphones.," in *Proc. INTERSPEECH*. ISCA, 2020, pp. 71–75.

[37] Siqi Zheng, Weilong Huang, Xianliang Wang, Hongbin Suo, Jinwei Feng, and Zhijie Yan, "A real-time speaker diarization system based on spatial spectrum," in *Proc. ICASSP*. IEEE, 2021, pp. 7208–7212.

[38] Chen Shen, Yi Liu, Wenzhi Fan, Bin Wang, Shixue Wen, Yao Tian, Jun Zhang, Jingsheng Yang, and Zejun Ma, "The volcspeech system for the icassp 2022 multi-channel multi-party meeting transcription challenge," in *Proc. ICASSP*. IEEE, 2022, pp. 9176–9180.

[39] Yihui Fu, Luyao Cheng, Shubo Lv, Yukai Jv, Yuxiang Kong, Zhuo Chen, Yanxin Hu, et al., "AISHELL-4: An open source dataset for speech enhancement, separation, recognition and speaker diarization in conference scenario," in *Proc. INTERSPEECH*. ISCA, 2021, pp. 3665–3669.

[40] Yanxin Hu, Yun Liu, Shubo Lv, Mengtao Xing, et al., "DCCRN: Deep complex convolution recurrent network for phase-aware speech enhancement," in *Proc. INTERSPEECH*. ISCA, 2020, pp. 2472–2476.

[41] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson-Enrique Yalta Soplin, Jahn Heymann, Matthew Wiesner, Nanxin Chen, et al., "ESPnet: End-to-End speech processing toolkit," in *Proc. INTERSPEECH*. ISCA, 2018, pp. 2207–2211.

[42] Shuaishuai Ye, Peiyao Wang, Shunfei Chen, Xinhui Hu, and Xinkang Xu, "The royalflush system of speech recognition for m2met challenge," in *Proc. ICASSP*. IEEE, 2022, pp. 9181–9185.