# Audio-visual End-to-end Multi-channel Speech Separation, Dereverberation and Recognition

Guinan Li, Jiajun Deng, Mengzhe Geng, Zengrui Jin, Tianzi Wang, Shujie Hu,
Mingyu Cui, Helen Meng, *Fellow, IEEE*, Xunying Liu, *Memeber, IEEE*

*Abstract*—Accurate recognition of cocktail party speech containing overlapping speakers, noise and reverberation remains a highly challenging task to date. Motivated by the invariance of visual modality to acoustic signal corruption, an audio-visual multi-channel speech separation, dereverberation and recognition approach featuring a full incorporation of visual information into all system components is proposed in this paper. The efficacy of the video input is consistently demonstrated in mask-based MVDR speech separation, DNN-WPE or spectral mapping (SpecM) based speech dereverberation front-end and Conformer ASR back-end. Audio-visual integrated front-end architectures performing speech separation and dereverberation in a pipelined or joint fashion via mask-based WPD are investigated. The error cost mismatch between the speech enhancement front-end and ASR back-end components is minimized by end-to-end jointly fine-tuning using either the ASR cost function alone, or its interpolation with the speech enhancement loss. Experiments were conducted on the mixture overlapped and reverberant speech data constructed using simulation or replay of the Oxford LRS2 dataset. The proposed audio-visual multi-channel speech separation, dereverberation and recognition systems consistently outperformed the comparable audio-only baseline by 9.1% and 6.2% absolute (41.7% and 36.0% relative) word error rate (WER) reductions. Consistent speech enhancement improvements were also obtained on PESQ, STOI and SRMR scores[1].

*Index Terms*—Audio-visual, Speech separation, Speech dereverberation, Speech recognition, End-to-end, Conformer

## I. INTRODUCTION

DESPITE the rapid progress of automatic speech recognition (ASR) in the past few decades, accurate recognition of cocktail party speech [1], [2] remains a highly challenging task to date. Its difficulty can be attributed to multiple sources of interference including overlapping speakers, background noise and room reverberation. These lead to a large mismatch between the resulting mixture speech and clean signals.

To this end, microphone arrays play a key role in state-of-the-art speech enhancement and recognition systems designed for cocktail party overlapped speech and far-field scenarios [3]–[5]. The required array beamforming techniques used to perform multi-channel signal integration are normally implemented as either time or frequency domain filters. These are represented by time domain delay and sum [6], frequency domain minimum variance distortionless response (MVDR) [7], [8] and generalized eigenvalue (GEV) [9] based multi-channel integration approaches. Earlier generations of mixed speech separation and recognition systems featuring conventional multi-channel array beamforming techniques typically used a pipelined system architecture. It contains separately constructed speech enhancement front-end modules designed to perform speech separation, dereverberation as well as denoising tasks, and speech recognition back-end components.

With the wider application of deep neural networks (DNNs) based speech technologies, microphone array beamforming techniques have also evolved into a rich variety of neural network based designs in recent few years. These include: a) neural time-frequency (TF) masking approaches [10]–[12] used to predict spectral mask labels for a reference channel that specify whether a particular TF spectrum point is dominated by the target speaker or interfering sources to facilitate speech separation; b) neural Filter and Sum approaches directly estimating the beamforming filter parameters in either time domain [13] or frequency domain [14] to produce the separated outputs; and c) mask-based MVDR [4], [15]–[19], and mask-based GEV [20], [21] approaches utilizing DNN estimated TF masks to compute target speaker and noise specific speech power spectral density (PSD) matrices and to obtain the beamforming filter parameters, while alleviating the need of explicit direction of arrival (DOA) estimation.

In many practical applications, reverberation presents a further challenge which can lead to severe speech recognition performance degradation [22], [23] when such systems are trained on anechoic and non-reverberant data. Classical solutions to the resulting dereverberation problem represented by, for example, weighted prediction error (WPE) [24], require the estimation of a time delayed linear filter. In recent years, there has been a similar trend of conventional speech dereverberation approaches [24]–[27] such as WPE evolving into their current DNN based variants. These include: a) the DNN-WPE [22], [23] method, which uses neural network estimated target signal PSD matrices in place of those traditionally obtained using maximum likelihood estimation trained complex value Gaussian Mixture Models [24] in the dereverberation filter estimation; and b) complex spectral masking [28], [29] and spectral mapping [30], [31] learning a transformation between reverberant and anechoic data.

End-to-end all neural microphone array based speech enhancement and recognition systems present a comprehensive and overarching solution to the cocktail party speech problem

Guinan Li, Jiajun Deng, Mengzhe Geng, Zengrui Jin, Tianzi Wang, Shujie Hu, Mingyu Cui are with the Chinese University of Hong Kong, China (email: {gnli, jjdeng, mzgeng, zrjin, twang, sjhu, mycui}@se.cuhk.edu.hk)

Helen Meng is with the Chinese University of Hong Kong, China (email: hmmeng@se.cuhk.edu.hk).

Xunying Liu is with the Chinese University of Hong Kong, China and the corresponding author (email: xyliu@se.cuhk.edu.hk).

[1] Enhanced audio examples for demonstration purposes are available in https://liguinan.github.io/AV-E2E-MC-ASR

by simultaneously performing speech separation, denoising and dereverberation. However, efforts on developing such systems are confronted by a number of key research challenges.

**1) Full incorporation of video modality:** Motivated by the bimodal nature of human speech perception and the invariance of visual information to extrinsic acoustic corruption, there has been a long history of developing audio-visual speech enhancement [32]–[49] and recognition [50]–[66] techniques. When processing the cocktail mixed speech, a holistic, consistent incorporation of visual information in all components of the entire system (speech separation, dereverberation and recognition) is preferred. In contrast, among existing researches, video information has mainly been partially incorporated into: a) the speech enhancement (separation and/or dereverberation) front-end [33]–[49] alone; or b) the speech recognition back-end [50]–[66] only. More recent works used video information in both the multi-channel speech separation and ASR [67], but not in speech dereverberation.

**2) Integration between speech separation and dereverberation modules:** Surface reflection of speech signals in reverberant environments distorts the DOA or TF-mask estimation for the target speaker. At the same time, interfering sound sources also impact the dereverberation filter estimation. Hence, a suitable form of integration between the speech separation and dereverberation techniques is required within the speech enhancement front-end sub-system. Possible integration solutions include: a) a pipelined architecture within which the speech separation and dereverberation components are sequentially connected in any order such as the previous researches in [21], [48], [68]; or b) a single architecture where both these two enhancement functions are implemented, for example, using weighted power minimization distortionless response (WPD) [69]–[71] and the related DNN TF-mask based WPD [72], [73] approaches. To date, such integration problem has only been investigated for audio-only speech enhancement [21], [69]–[77], but has not been studied for audio-visual speech separation and dereverberation.

**3) Joint optimization of audio-visual speech enhancement front-end and recognition back-end:** Conventional non-DNN based speech enhancement front-end models are often separately constructed and cannot be easily integrated with the ASR back-end. The wide application of deep learning approaches for speech enhancement and recognition components allows them to be more tightly integrated and consistently optimized in an end-to-end manner. An improved trade-off between the speech enhancement front-end loss function and ASR accuracy can then be obtained, for example, using multi-task learning [67], [78], [79]. To date, such joint speech enhancement front-end and ASR back-end optimization has been only conducted among: a) audio-only speech enhancement and recognition systems using no video input [19], [23], [72], [78], [80]–[82]; or b) audio-visual speech separation and recognition tasks only while not considering speech dereverberation [67], [79]. Hence, there is a pressing need to derive suitable joint optimization methods for a complete audio-visual multi-channel speech separation, dereverberation and recognition system.

In order to address the above issues, an audio-visual multi-channel speech separation, dereverberation and recognition approach featuring a full incorporation of visual information into all three components of the entire system is proposed in this paper. The efficacy of the video input is consistently demonstrated when being used in the mask-based MVDR speech separation, DNN-WPE or spectral mapping (SpecM) based speech dereverberation front-end and Conformer encoder-decoder based ASR back-end components. Both the pipelined integration methods using either a) a serial connection of the audio-visual speech separation component with the following dereverberation module; or b) audio-visual speech dereverberation followed by separation; and c) joint speech separation and dereverberation via audio-visual mask-based WPD are investigated. In order to reduce the error cost mismatch between the speech enhancement front-end and ASR back-end components, they are jointly fine-tuned using either only the Conformer ASR cost function (CTC plus Attention) [83], or the ASR cost function interpolated with the speech enhancement loss based on mean square error (MSE) and scale-invariant signal to noise ratio (SISNR).

Experiments conducted on the mixture overlapped and reverberant speech data constructed using either simulation or replay of the benchmark Oxford LRS2 dataset [84] suggest:

1) The proposed audio-visual multi-channel speech separation, dereverberation and recognition systems consistently outperformed the comparable audio-only baseline systems by **9.1% and 6.2% absolute (41.7% and 36.0% relative)** word error rate (WER) reductions on the LRS2 simulated and replayed evaluation datasets, respectively. Consistent improvements of perceptual evaluation of speech quality (PESQ) [85], short-time objective intelligibility (STOI) [86] and speech to reverberation modulation energy ratio (SRMR) [87] scores were also obtained.

2) In particular, when compared with audio-only dereverberation, incorporating visual information into the DNN-WPE or SpecM based dereverberation module produced consistent improvements of PESQ, STOI and SRMR scores and a statistically significant[2] WER reduction by up to **1.9% absolute (5.9% relative)**, irrespective of the form of integration between speech separation and dereverberation components.

3) Among different architectures to integrate the speech separation and dereverberation components within the front-end, a pipelined, full audio-visual configuration performing DNN-WPE based speech dereverberation followed by mask-based MVDR speech separation using video input in both stages produced the best overall speech enhancement and recognition performance.

4) Consistent WER reductions and improvements on speech enhancement metric scores were also obtained after joint fine-tuning the entire audio-visual speech separation, dereverberation and recognition system in a fully end-to-end manner.

The main contributions of this paper are summarized below:

1) To the best of our knowledge, this paper presents the first use of a complete audio-visual multi-channel speech separation, dereverberation and recognition system architecture featuring a full incorporation of visual information into all

---

[2]Matched pairs sentence-segment word error (MAPSSWE) based statistical significance test [88] was performed at a significance level $\alpha$=0.05.

Fig. 1. Audio-visual multi-channel **speech separation** using mask-based MVDR approach (a), and **joint speech separation & dereverberation module** using mask-based WPD in (b). Both use the same audio-visual embeddings (left part of the figure) for their complex masks estimation. $Y_r(t, f) \in \mathbb{C}$ is the $r$-th channel's complex spectrum of mixture speech among $R$ microphone channels. $\mathbf{V}(t)$ and $\mathbf{A}(t)$ denote the audio and visual embeddings at frame index $t$, respectively. The internal structural details of the TCN block and Visual Conv1DBlock are shown in Fig. 2. The MVDR filter $\mathbf{w}_{\text{MVDR}}(f) \in \mathbb{C}^R$ is estimated using the target speech and noise PSD matrices $\mathbf{\Phi}_x(f) \in \mathbb{C}^{R \times R}$ and $\mathbf{\Phi}_n(f) \in \mathbb{C}^{R \times R}$ with their respective complex TF masks $M_{\text{MVDR}}^x(t, f) \in \mathbb{C}$ and $M_{\text{MVDR}}^n(t, f) \in \mathbb{C}$. The WPD filter $\tilde{\mathbf{w}}_{\text{WPD}}(f) \in \mathbb{C}^{(L+1)R}$ is estimated using the target speaker and power normalized spatial-temporal PSD matrices $\mathbf{\Phi}_{\tilde{x}}(f) \in \mathbb{C}^{(L+1)R \times (L+1)R}$ and $\mathbf{\Phi}_{\tilde{y}}(f) \in \mathbb{C}^{(L+1)R \times (L+1)R}$ with their respective complex TF masks $M_{\text{WPD}}^{\tilde{x}}(t, f) \in \mathbb{C}$ and $M_{\text{WPD}}^\lambda(t, f) \in \mathbb{C}$. Re($\cdot$) and Im($\cdot$) denote the real and imaginary parts operators. $D$ is the prediction delay parameter and $L$ is the number of filter taps.

three stages. In contrast, prior researches incorporate visual modality in either only the speech enhancement front-end [33]–[49], ASR back-end [50]–[66], or both the multi-channel speech separation and recognition stages [67] but excluding the dereverberation component.

2) This paper presents a more complete investigation of the advantages of audio-visual dereverberation approaches versus audio-only dereverberation methods based on DNN-WPE and SpecM. In contrast, similar prior studies [48] were conducted only in the context of SpecM based dereverberation.

3) To the best of our knowledge, this is the first work that systematically investigates the suitable form of integration between the full audio-visual speech separation and dereverberation modules within the speech enhancement front-end. In contrast, similar studies in previous researches were only conducted for audio-only speech enhancement [72].

4) This paper presents the first research to demonstrate that performing an end-to-end joint optimization is useful for training a complete audio-visual multi-channel speech separation, dereverberation and recognition system. In contrast, related prior studies were conducted only in the context of audio-only speech enhancement and recognition [72].

We hope these findings above will provide valuable insights for the practical development of state-of-the-art audio-visual speech separation, dereverberation and recognition systems for cocktail party and far-field scenarios.

The rest of the paper is organized as follows. Audio-visual multi-channel speech separation is reviewed in Section II. Section III presents audio-visual multi-channel speech dereverberation. Integrated audio-visual speech separation and dereverberation approaches are proposed in Section IV. Section V presents the audio-visual Conformer ASR back-end component and its joint fine-tuning with the speech enhancement front-end. Experimental data setup and results are presented in Section VI and VII, respectively. Section VIII draws the conclusion and discusses future research directions.

## II. AUDIO-VISUAL MULTI-CHANNEL SPEECH SEPARATION

In this section, the multi-channel far-field speech signal model is reviewed first, before the introduction of the audio-visual multi-channel mask-based MVDR approach for speech separation is presented.

### A. Multi-channel Far-field Signal Model

In the far-field scenarios, the short-time Fourier transform (STFT) spectrum of the received multi-channel speech signal $\mathbf{y}(t, f) \in \mathbb{C}^R$ recorded by a microphone array consisting of $R$ channels can be modeled as:

$$\mathbf{y}(t, f) = \mathbf{x}(t, f) + \mathbf{n}(t, f) = \mathbf{g}(f)S(t, f) + \mathbf{n}(t, f), \quad (1)$$

where $t$ and $f$ denote the indices of time and frequency bins, respectively. $\mathbf{x}(t, f) \in \mathbb{C}^R$ is a complex vector containing the clean speech signals received by the array channels. $\mathbf{n}(t, f) \in \mathbb{C}^R$ represents either the interfering speaker's speech or additive background noise alone, or a combination of both. $\mathbf{g}(f) \in \mathbb{C}^R$ denotes the array steering vector and $S(t, f)$ is the STFT spectrum of the target speaker's clean speech.

### B. Mask-based MVDR

Classic acoustic beamforming approaches [7]–[9] are designed to capture the speech from the target speaker's direction while attenuating the interfering sounds coming from other locations. This is realized by setting, or "steering", the beamforming filter parameters to the target direction. Taking the MVDR beamformer as an example, a linear filter $\mathbf{w}_{\text{MVDR}}(f) \in \mathbb{C}^R$ is applied to the multi-channel mixture speech spectrum $\mathbf{y}(t, f)$ to produce the filtered output $\hat{S}_{\text{MVDR}}(t, f)$ as:

$$\hat{S}_{\text{MVDR}}(t, f) = \mathbf{w}_{\text{MVDR}}(f)^H \mathbf{y}(t, f), \quad (2)$$

$$= \underbrace{\mathbf{w}_{\text{MVDR}}(f)^H \mathbf{x}(t, f)}_{\text{target speech component}} + \underbrace{\mathbf{w}_{\text{MVDR}}(f)^H \mathbf{n}(t, f)}_{\text{residual noise}}, \quad (3)$$

where $(\cdot)^H$ denotes the conjugate transpose operator.

The MVDR beamformer is designed to minimize the residual noise output while imposing a distortionless constraint on the target speech [7], which can be formulated as

$$\min_{\mathbf{w}_{\mathrm{MVDR}}(f)} \sum_t \left| \mathbf{w}_{\mathrm{MVDR}}(f)^H \mathbf{n}(t,f) \right|^2, \tag{4}$$

$$\text{subject to}: \sum_t \left| (\mathbf{u}_r - \mathbf{w}_{\mathrm{MVDR}}(f))^H \mathbf{x}(t,f) \right|^2 = 0, \tag{5}$$

where $\mathbf{u}_r = [0, 0, \ldots, 1, \ldots, 0]^T \in \mathbb{R}^R$ is a one-hot reference vector where its $r$-th component equals to one. $(\cdot)^T$ denotes the transpose operator. Without loss of generality, we select the first channel, i.e., $r = 1$ as the reference channel among the $R$ channels throughout this paper.

The distortionless constraint in the above optimization problem is equivalent to $\mathbf{w}_{\mathrm{MVDR}}(f)^H \mathbf{g}(f) = 1$, which can be interpreted as maintaining the energy along the target direction. The MVDR beamforming filter is estimated as

$$\mathbf{w}_{\mathrm{MVDR}}(f) = \frac{\boldsymbol{\Phi}_n(f)^{-1} \mathbf{g}(f)}{\mathbf{g}(f)^H \boldsymbol{\Phi}_n(f)^{-1} \mathbf{g}(f)} = \frac{\boldsymbol{\Phi}_n(f)^{-1} \boldsymbol{\Phi}_x(f)}{\mathrm{tr}\left( \boldsymbol{\Phi}_n(f)^{-1} \boldsymbol{\Phi}_x(f) \right)} \mathbf{u}_r, \tag{6}$$

where the target speaker and noise specific power spectral density (PSD) matrices

$$\boldsymbol{\Phi}_x(f) = \frac{\sum_t \left( M_{\mathrm{MVDR}}^x(t,f)\mathbf{y}(t,f) \right) \left( M_{\mathrm{MVDR}}^x(t,f)\mathbf{y}(t,f) \right)^H}{\sum_t M_{\mathrm{MVDR}}^x(t,f) \left( M_{\mathrm{MVDR}}^x(t,f) \right)^*}, \tag{7}$$

$$\boldsymbol{\Phi}_n(f) = \frac{\sum_t \left( M_{\mathrm{MVDR}}^n(t,f)\mathbf{y}(t,f) \right) \left( M_{\mathrm{MVDR}}^n(t,f)\mathbf{y}(t,f) \right)^H}{\sum_t M_{\mathrm{MVDR}}^n(t,f) \left( M_{\mathrm{MVDR}}^n(t,f) \right)^*}, \tag{8}$$

are computed using DNN predicted complex TF masks $M_{\mathrm{MVDR}}^x(t,f) \in \mathbb{C}$ and $M_{\mathrm{MVDR}}^n(t,f) \in \mathbb{C}$ [19], [67]. $\mathrm{tr}(\cdot)$ denotes the trace operator. $(\cdot)^*$ is complex conjugate operator.

### C. Audio Modality

As is illustrated in the top left corner of Fig. 1, three types of audio features including the complex STFT spectrum of all the microphone array channels, the inter-microphone phase differences (IPDs) [15] and location-guided angle feature (AF) [89] are adopted as the audio inputs. IPDs features are used to capture the relative phase difference between different microphone channels and provide additional spatial cues for mask-based multi-channel speech separation. Angle features that are based on the approximated DOA of the target speaker[3] are also incorporated to provide further spatial filtering constraints. In this work, the approximated DOA of the target speaker is obtained by tracking the speaker's face from a $180°$ wide-angle camera (Fig. 1, bottom left corner).

Following prior researches on audio-visual multi-channel speech separation [67], [68], the temporal convolutional network architecture (TCN) [90], which uses a long reception field to capture more sufficient contextual information, is used in our separation system. As shown in the left of Fig. 2, each TCN block is stacked by 8 Dilated 1-D ConvBlock with exponentially increased dilation factors $2^0, 2^1, \ldots, 2^7$. As shown in the top left corner of Fig. 1, the log-power spectrum (LPS) features of the reference microphone channel are

concatenated with the IPDs and AF features before being fed into a single TCN module based Audio Block to compute the audio embeddings $\mathbf{A} \in \mathbb{R}^{F_a \times T_a}$, where $F_a$ is the dimension of audio embeddings and $T_a$ is the number of audio frames.



Fig. 2. Illustration of the architectures of: (a) the temporal convolutional network (TCN) Block. Each dilated 1-D ConvBlock consists of a $1 \times 1$ convolutional layer, a depth-wise separable convolution layer (D-Conv) [91], with PReLU [92] activation function and batch normalization added between each two convolution layers. Skip connections are added in each dilated 1-D ConvBlock; and (b) Visual Conv1DBlock which consists of a PReLU [92] activation function, batch normalization, a depth-wise separable convolution layer (D-Conv) [91] and a $1 \times 1$ convolutional layer with skip connection.

### D. Visual Modality

The lip region of a target speaker obtained via face tracking is fed into a LipNet [93] which consists of a 3D convolutional layer (Fig. 1, bottom left, in pink) and an 18-layer ResNet [94] (Fig. 1, bottom left, in light turquoise), to extract the visual features from the target speaker's lip movements. Before fusing the visual features with the audio embeddings to improve the TF masks estimation, the visual features are firstly fed into the linear layer followed by the Visual Block containing five Visual Conv1DBlocks (Fig. 1, bottom, in light brown, the detailed network architecture is illustrated in the right of Fig. 2), and then the output of Visual Block is up-sampled to be time synchronised with the audio frames via linear interpolation to compute the visual embeddings $\mathbf{V} \in \mathbb{R}^{F_v \times T_a}$, where $F_v$ is the dimension of visual embeddings. In this work, the LipNet model is pretrained on the lipreading task as described in [93].

### E. Modality Fusion

In order to effectively integrate the audio and visual embeddings, a factorized attention-based modality fusion method [67], [68] is utilized in the audio-visual speech separation module. As shown in Fig. 1 (middle up), the acoustic embeddings at frame index $t$ denoted by $\mathbf{A}(t)$ are first factorized into $K$ acoustic subspace vectors $[\mathbf{e}_1^a(t), \mathbf{e}_2^a(t), \ldots, \mathbf{e}_K^a(t)]$ by a series of parallel linear transformation $\mathbf{P}_k^a \in \mathbb{R}^{F_a \times F_a}$. The visual embeddings at frame index $t$ named by $\mathbf{V}(t)$ is mapped into a $K$ dimensional vector $\mathbf{e}^v(t) = [e_1^v(t), e_2^v(t), \ldots, e_K^v(t)]^T$ by projection matrix $\mathbf{P}^v \in \mathbb{R}^{K \times F_v}$ as

$$[\mathbf{e}_1^a(t), \mathbf{e}_2^a(t), \ldots, \mathbf{e}_K^a(t)] = [\mathbf{P}_1^a, \mathbf{P}_2^a, \ldots, \mathbf{P}_K^a] \mathbf{A}(t), \tag{9}$$

$$\mathbf{e}^v(t) = \mathrm{Softmax}\left( \mathbf{P}^v \mathbf{V}(t) \right), \tag{10}$$

Then the fused audio-visual embeddings $\mathbf{AV}(t) \in \mathbb{R}^{F_a}$ are

$$\mathbf{AV}(t) = \boldsymbol{\sigma} \left( \sum_{k=1}^{K} e_k^v(t) \mathbf{e}_k^a(t) \right), \tag{11}$$

where $\boldsymbol{\sigma}(\cdot)$ is the sigmoid function.

---

[3]The target speaker is located using a 180-degree wide-angle camera to track the speaker's face. The camera approximated DOA of target speaker is only used in AF features.

Fig. 3. Illustration of audio-visual multi-channel speech dereverberation networks based on the **(a) DNN-WPE** or **(b) SpecM** approaches of Sections III-B and III-C respectively. $X_r(t, f) \in \mathbb{C}$ is the $r$-th channel's complex spectrum of reverberant speech among R microphone channels. $\mathbf{V}(t)$ and $\mathbf{A}(t)$ denote the audio and visual embeddings at frame index $t$, in common with Fig.1. During WPE filter estimation, the signal variance $\lambda(t, f)$ is obtained using DNN predicted TF complex mask $M_{\text{WPE}}(t, f) \in \mathbb{C}$. $\mathbf{x}(t, f) \in \mathbb{C}^R$ is the input multi-channel reverberant speech signal. $D$ denotes the prediction delay parameter and $L$ is the number of filter taps. $M_{\text{SpecM}}(t, f) \in \mathbb{C}$ denotes the complex TF mask for SpecM based dereverberation.

The above audio-visual embeddings are fed into both the Target Speech Block and Noise Block (Fig. 1, center), before their respective outputs being further fed into the corresponding linear layers (Fig. 1, top right, yellow blocks) to estimate the complex TF masks $M_{\text{MVDR}}^x(t, f) \in \mathbb{C}$ and $M_{\text{MVDR}}^n(t, f) \in \mathbb{C}$ required by the target speech and noise PSD matrices in Eqns. (7) and (8) for MVDR filter estimation. After MVDR filtering, the separated target speech spectrum is inverse STFT (iSTFT) transformed to produce the corresponding waveform.

### F. Separation Network Training Cost Function

Following the prior researches [36], [48], [67], [68], the mask-MVDR based multi-channel speech separation network is trained to maximize the SISNR metric, unless further joint fine-tuning with the back-end ASR error loss later presented in Section V is performed.

## III. AUDIO-VISUAL MULTI-CHANNEL SPEECH DEREVERBERATION

In this section, the multi-channel far-field signal model is reformulated with additional reverberation. Audio-visual multi-channel speech dereverberation approaches based on audio-visual DNN-WPE and SpecM are then proposed. The incorporation of the video features and its fusion with audio modality in both methods are also presented.

### A. Multi-channel Far-field Signal Model with Reverberation

In reverberant conditions, the target speech signal $\mathbf{x}(t, f)$ of Eqn. (1) is further decomposed into two parts. The first part consists of the direct signal and early reflections, referred to as the desired signal $\mathbf{d}(t, f) \in \mathbb{C}^R$, while the other contains the late reverberation $\mathbf{r}(t, f) \in \mathbb{C}^R$. This is given by

$$\mathbf{x}(t,f) = \underbrace{\sum_{\tau=0}^{D-1} \mathbf{a}(\tau, f) S(t-\tau, f)}_{\mathbf{d}(t,f)} + \underbrace{\sum_{\tau=D}^{D+L-1} \mathbf{a}(\tau, f) S(t-\tau, f)}_{\mathbf{r}(t,f)} \quad (12)$$

where $D$ denotes the prediction delay parameter and $L$ is the number of filter taps. $\mathbf{a}(\tau, f) \in \mathbb{C}^R$ is the room reverberant transfer function from a given speaker to all microphones for $\tau \in \{0, 1, \ldots, D+L-1\}$. The dereverberation process requires the desired signal $\mathbf{d}(t, f)$ to be preserved to enhance speech intelligibility and improve ASR performance, while the late reverberation $\mathbf{r}(t, f)$ to be eliminated [24].

### B. DNN-WPE Based Dereverberation

In conventional WPE [24], the dereverberated signal $\hat{\mathbf{d}}(t, f)$ can be obtained by applying the WPE filter $\mathbf{W}_{\text{WPE}}(f) \in \mathbb{C}^{LR \times R}$ to the reverberant multi-channel signal as follows:

$$\hat{\mathbf{d}}(t, f) = \mathbf{x}(t, f) - \mathbf{W}_{\text{WPE}}(f)^H \tilde{\mathbf{x}}(t - D, f), \quad (13)$$

where $\tilde{\mathbf{x}}(t-D, f) = \left[\mathbf{x}(t-D, f)^T, \ldots, \mathbf{x}(t-D-L+1, f)^T\right]^T \in \mathbb{C}^{LR}$ is the time-delayed reverberant speech spectrum vector.

The required WPE filter coefficients are traditionally estimated using maximum likelihood estimation [24]. It is assumed that the desired signal at each microphone follows a time-varying complex Gaussian distribution with a mean of zero and a time-varying variance $\lambda(t, f)$, which corresponds to the power of the desired signal. Minimizing the average power of the frame prediction errors weighted by $\lambda^{-1}(t, f)$,

$$\min_{\{\mathbf{W}_{\text{WPE}}(f), \lambda(t,f)\}} \sum_t \frac{\left\|\mathbf{x}(t, f) - \mathbf{W}_{\text{WPE}}(f)^H \tilde{\mathbf{x}}(t-D, f)\right\|_2^2}{\lambda(t, f)}. \quad (14)$$

leads to alternating updates between the WPE filter parameters,

$$\mathbf{W}_{\text{WPE}}(f) = \left(\sum_t \frac{\tilde{\mathbf{x}}(t-D, f)\tilde{\mathbf{x}}(t-D, f)^H}{\lambda(t, f)}\right)^{-1} \left(\sum_t \frac{\tilde{\mathbf{x}}(t-D, f)\mathbf{x}(t, f)^H}{\lambda(t, f)}\right) \quad (15)$$

and the residual signal power given the current WPE filter

$$\lambda(t, f) = \frac{1}{R} \left\|\hat{\mathbf{d}}(t, f)\right\|_2^2, \quad (16)$$

where $\|\cdot\|_2$ denotes the Euclidean norm. The above alternating estimation procedure iterates until convergence.

Recent deep neural network extension to WPE led to the DNN-WPE approach [22], where the filtered signal power $\lambda(t, f)$ is estimated using DNN (e.g. LSTM [22]) predicted TF complex mask[4] $M_{\text{WPE}}(t, f) \in \mathbb{C}$. This is given by

$$\lambda(t, f) = \frac{1}{R} \|M_{\text{WPE}}(t, f)\mathbf{x}(t, f)\|_2^2, \quad (17)$$

An example of DNN-WPE based dereverberation is shown in Fig. 3 (top right, in light blue).

[4]Alternatively using channel dependent predicted mask $M_{\text{WPE}}^r(t, f)$ produced comparable performance in practice while increasing the system training time approximately by a factor of 5, and therefore not considered.

### C. SpecM Based Dereverberation

In addition to DNN-WPE based dereverberation, SpecM based dereverberation is also leveraged in this work. A neural network based TF spectral transformation between the input reverberant and desired anechoic speech spectrum is learned as follows:

$$\hat{\mathbf{d}}(t,f) = W_{\text{SpecM}}(t,f)\mathbf{x}(t,f) = M_{\text{SpecM}}(t,f)\mathbf{x}(t,f), \quad (18)$$

where $W_{\text{SpecM}}(t,f) \in \mathbb{C}$ denotes the SpecM filter and $M_{\text{SpecM}}(t,f) \in \mathbb{C}$ is the estimated complex TF-mask for SpecM based dereverberation.

An example of SpecM based speech dereverberation is shown in Fig. 3 (bottom right, in light yellow). Compared with DNN-WPE, although the SpecM based dereverberation approach can provide perceptually enhanced sounds, it has been reported that the artifacts resulting from deterministic spectral masking introduced a negative impact on downstream speech recognition system performance [3], [15], [16].

### D. Audio-visual Speech Dereverberation

The audio and video embeddings previously used in the mask-based MVDR speech separation network of Section II and Fig. 1 are concatenated[5] before being fed into an AV Fusion Block consisting of three TCN modules to produce the integrated audio-visual embeddings (Fig. 3, left).

These audio-visual embeddings are then forwarded into linear layers (Fig. 3, right, yellow blocks) to estimate the complex TF masks of the desired speech for either DNN-WPE (Fig. 3, top right, light blue) or SpecM (Fig. 3, bottom right, light yellow) based dereverberation filter estimation. In this work, the dereverberation network is trained in both cases using the MSE loss computed between the filtered and ground-truth anechoic speech spectrum [22], [48], [68].

### IV. AUDIO-VISUAL SEPARATION AND DEREVERBERATION

In this section, three integrated audio-visual speech separation and dereverberation architectures are proposed. These include: a) a serial pipelined connection of the audio-visual speech separation component with the following dereverberation module; or b) conversely audio-visual speech dereverberation followed by separation; and c) joint speech separation & dereverberation using audio-visual mask-based WPD.

### A. Audio-visual Speech Separation-Dereverberation

In the audio-visual speech separation-dereverberation architecture, the multi-channel mixture speech spectra $\mathbf{y}(t,f) \in \mathbb{C}^R$ as well as the extracted visual features and the camera captured target speaker's DOA from the Visual Front-end module (e.g. Fig. 1, bottom left corner, in light green) are first fed into the MVDR separation module as shown in Fig. 1(a) to produce single-channel outputs, $\hat{S}_{\text{MVDR}}(t,f)$, before being connected to the dereverberation module based on DNN-WPE or SpecM as shown in Fig. 3 to obtain the final enhanced speech $\hat{d}_{\text{MVDR-WPE}}(t,f) \in \mathbb{C}$ or $\hat{d}_{\text{MVDR-SpecM}}(t,f) \in \mathbb{C}$, respectively.

---

[5]Alternative audio-visual modality fusion methods, e.g. using the factorized attention based fusion mechanism of Section II-E for speech separation, led to performance degradation in practice and therefore not considered.

When DNN-WPE based dereverberation is used, this is computed in a two stage, pipelined manner as

$$\hat{S}_{\text{MVDR}}(t,f) = \mathbf{w}_{\text{MVDR}}(f)^H \mathbf{y}(t,f), \quad (19)$$

$$\hat{d}_{\text{MVDR-WPE}}(t,f) = \hat{S}_{\text{MVDR}}(t,f) - \mathbf{W}_{\text{WPE}}(f)^H \hat{\mathbf{s}}_{\text{MVDR}}(t-D,f), \quad (20)$$

where

$$\hat{\mathbf{s}}_{\text{MVDR}}(t-D,f) = \left[\hat{S}_{\text{MVDR}}(t-D,f), \ldots, \hat{S}_{\text{MVDR}}(t-D-L+1,f)\right]^T$$

denotes the enhanced single-channel output of the MVDR beamformer from the past $L$ frames and $\hat{\mathbf{s}}_{\text{MVDR}}(t-D,f) \in \mathbb{C}^L$. Here, $\mathbf{W}_{\text{WPE}}(f) \in \mathbb{C}^L$ represents the single-channel WPE filter. $L$ is the number of filter taps and $D$ denotes the prediction delay parameter in WPE.

When SpecM based dereverberation is used, the final enhanced single-channel speech spectrum is computed as

$$\hat{S}_{\text{MVDR}}(t,f) = \mathbf{w}_{\text{MVDR}}(f)^H \mathbf{y}(t,f), \quad (21)$$

$$\hat{d}_{\text{MVDR-SpecM}}(t,f) = W_{\text{SpecM}}(t,f)\hat{S}_{\text{MVDR}}(t,f). \quad (22)$$

### B. Audio-visual Speech Dereverberation-Separation

In contrast to the above, connecting the speech dereverberation and separation modules in a reverse order leads to the audio-visual speech dereverberation-separation architecture. The sequence of filtering operations of this architecture is performed as follows:

When using DNN-WPE based dereverberation, the dereverberated multi-channel output $\hat{\mathbf{d}}_{\text{WPE}}(t,f)$ is first produced, before being fed into the MVDR separation filter to produce the final single-channel speech spectrum $\hat{S}_{\text{WPE-MVDR}}(t,f)$ as

$$\hat{\mathbf{d}}_{\text{WPE}}(t,f) = \mathbf{y}(t,f) - \mathbf{W}_{\text{WPE}}(f)^H \tilde{\mathbf{y}}(t-D,f), \quad (23)$$

$$\hat{S}_{\text{WPE-MVDR}}(t,f) = \mathbf{w}_{\text{MVDR}}(f)^H \hat{\mathbf{d}}_{\text{WPE}}(t,f), \quad (24)$$

where $\tilde{\mathbf{y}}(t-D,f) = \left[\mathbf{y}(t-D,f)^T, \ldots, \mathbf{y}(t-D-L+1,f)^T\right]^T \in \mathbb{C}^{LR}$ denotes the stacked vector representation of the input multi-channel mixture speech signal.

When using SpecM based dereverberation, the above can be expressed as

$$\hat{\mathbf{d}}_{\text{SpecM}}(t,f) = W_{\text{SpecM}}(t,f)\mathbf{y}(t,f), \quad (25)$$

$$\hat{S}_{\text{SpecM-MVDR}}(t,f) = \mathbf{w}_{\text{MVDR}}(f)^H \hat{\mathbf{d}}_{\text{SpecM}}(t,f). \quad (26)$$

### C. Audio-visual Joint Speech Separation & Dereverberation

Combining the multi-channel speech separation and dereverberation functions into a single convolutional filter leads to a joint speech separation and dereverberation architecture, for example, based on WPD [69]–[71] and their DNN predicted mask-based variants [72].

When producing the final enhanced speech spectrum, a single WPD filter $\tilde{\mathbf{w}}_{\text{WPD}}(f) \in \mathbb{C}^{(L+1)R}$ is applied to the time-delayed multi-channel mixed speech vector stacked by $\mathbf{y}(t,f) \in \mathbb{C}^R$ and $\tilde{\mathbf{y}}(t-D,f)^T \in \mathbb{C}^{LR}$ as follows:

$$\hat{d}(t,f) = \tilde{\mathbf{w}}_{\text{WPD}}(f)^H \left[\mathbf{y}(t,f)^T, \tilde{\mathbf{y}}(t-D,f)^T\right]^T, \quad (27)$$

The WPD beamformer is trained to minimize the average weighted power of the filtered signal while satisfying

Fig. 4. Illustration of an end-to-end audio-visual multi-channel speech separation, dereverberation and recognition system, which integrates the Speech Enhancement Front-end, Visual Front-end, Feature Extraction and Conformer ASR Back-end components.

an orthogonal constraint for channel synchronization without distorting the target speech. This is given by

$$\min_{\tilde{\mathbf{w}}_{\text{WPD}}(f)} \sum_t \frac{\left| \tilde{\mathbf{w}}_{\text{WPD}}(f)^H \left[ \mathbf{y}(t,f)^T, \tilde{\mathbf{y}}(t-D,f)^T \right]^T \right|^2}{\lambda(t,f)}, \quad (28)$$

$$\text{subject to} \;: \tilde{\mathbf{w}}_{\text{WPD}}(f)^H \tilde{\mathbf{g}}(f) = 1. \quad (29)$$

where the signal variance is averaged across $R$ channels as

$$\lambda(t,f) = \frac{1}{R} \sum_{r=1}^{R} \left| M_{\text{WPD}}^{\lambda}(t,f) Y_r(t,f) \right|^2,$$

is estimated using DNN predicted TF complex mask of the desired signal $M_{\text{WPD}}^{\lambda}(t,f) \in \mathbb{C}$. $Y_r(t,f)$ represents the $r$-th component of the multi-channel mixture speech signal $\mathbf{y}(t,f)$. $\tilde{\mathbf{g}}(f) = \left[ \mathbf{g}(f)^T, \mathbf{0}, \ldots, \mathbf{0} \right]^T \in \mathbb{C}^{(L+1)R}$ is the padded steering vector which is composed of a steering vector $\mathbf{g}(f) \in \mathbb{C}^R$ and the others $\mathbf{0} \in \mathbb{C}^R$ vectors. It can be shown that the solution of the above WPD convolutional beamformer is:

$$\tilde{\mathbf{w}}_{\text{WPD}}(f) = \frac{\mathbf{\Phi}_{\tilde{y}}(f)^{-1} \tilde{\mathbf{g}}(f)}{\tilde{\mathbf{g}}(f)^H \mathbf{\Phi}_{\tilde{y}}(f)^{-1} \tilde{\mathbf{g}}(f)} = \frac{\mathbf{\Phi}_{\tilde{y}}(f)^{-1} \mathbf{\Phi}_{\tilde{x}}(f)}{\text{tr} \left( \mathbf{\Phi}_{\tilde{y}}(f)^{-1} \mathbf{\Phi}_{\tilde{x}}(f) \right)} \tilde{\mathbf{u}}_r, \quad (30)$$

where the target speaker and power normalized spatial-temporal PSD matrices are

$$\mathbf{\Phi}_{\tilde{x}}(f) = \frac{\sum_t \left( M_{\text{WPD}}^{\tilde{x}}(t,f) \tilde{\mathbf{y}}(t,f) \right) \left( M_{\text{WPD}}^{\tilde{x}}(t,f) \tilde{\mathbf{y}}(t,f) \right)^H}{\sum_t M_{\text{WPD}}^{\tilde{x}}(t,f) \left( M_{\text{WPD}}^{\tilde{x}}(t,f) \right)^*}, \quad (31)$$

$$\mathbf{\Phi}_{\tilde{y}}(f) = \sum_t \frac{\tilde{\mathbf{y}}(t,f) \tilde{\mathbf{y}}(t,f)^H}{\lambda(t,f)}, \quad (32)$$

and $\tilde{\mathbf{y}}(t,f) = \left[ \mathbf{y}(t,f)^T, \tilde{\mathbf{y}}(t-D,f)^T \right]^T \in \mathbb{C}^{(L+1)R}$. $\tilde{\mathbf{u}}_r = \left[ \mathbf{u}_r, \mathbf{0}, \ldots, \mathbf{0} \right]^T$ is the padded reference vector. $M_{\text{WPD}}^{\tilde{x}}(t,f) \in \mathbb{C}$ denotes the complex TF mask of target speech.

An example of mask-based WPD is illustrated in Fig. 1(b) (bottom right, in light blue). The same audio-visual embeddings that are used in mask-based MVDR separation module (Fig. 1, top right, light yellow) are now fed into three TCN based Target Speech Block and Time-varying Power Block for WPD filtering. Their respective outputs are then fed into the separate linear layers to estimate the complex TF masks $M_{\text{WPD}}^{\tilde{x}}(t,f) \in \mathbb{C}$ and $M_{\text{WPD}}^{\lambda}(t,f) \in \mathbb{C}$ required for the computation of the two spatial-temporal PSD matrices and finally the WPD filter parameters. The entire mask-based WPD network is trained using an equally weighted interpolation between the SISNR and MSE losses to perform joint speech separation & dereverberation.

## V. AUDIO-VISUAL MULTI-CHANNEL SPEECH RECOGNITION

In this section, the Conformer-based audio-visual speech recognition back-end and its further integration with the speech enhancement front-end are introduced.

### A. Audio-visual Conformer Speech Recognition Back-end

As shown in Fig. 4 (bottom left), the enhanced speech waveform produced by the speech separation and dereverberation front-ends of Sections II, III and IV is fed through a STFT transform before log Mel-filterbank (Mel-FBK) audio features are calculated. As is also shown in Fig. 4 (top left), the visual features extracted from the Visual Front-end are forwarded into a linear layer before being up-sampled to be time synchronised with the Mel-FBK audio frames. Finally, the audio and visual features are concatenated and fed into the ASR back-end.

The Conformer ASR back-end [95], [96] comprises a Conformer encoder and a Transformer decoder. The Conformer encoder has one convolutional subsampling module, and a linear layer with dropout operation followed by stacked encoder blocks. The internal components of each Conformer encoder block include: a position-wise feed-forward network module, a multi-head self-attention module, a convolution module, and a final position-wise feed-forward network module at the end. All the encoder blocks additionally undergo layer normalization and residual connections. Fig. 4 (right) shows an example of a Conformer ASR system, where the backbone model architecture is in the grey colored part (Fig. 4, bottom right). The detailed encoder block compositions are in the blue colored part (Fig. 4, top right). The following multi-task criterion interpolation between the CTC and attention error costs [83] is utilized in Conformer model training,

$$\mathcal{L}_{\text{ASR}} = (1 - \beta) \mathcal{L}_{att} + \beta \mathcal{L}_{ctc}, \quad (33)$$

where $\beta \in [0, 1]$ is a tunable hyper-parameter and empirically set as 0.3 for training and 0.4 for recognition in this paper.

### B. Integration of Speech Enhancement and Recognition

Traditionally, the speech enhancement front-end and recognition back-end components are optimized separately and used in a pipelined manner [15], [16], [21], [97], [98]. However, two issues arise with this pipelined approach: **1)** the learning cost function mismatch between speech enhancement front-end and recognition back-end components is not addressed; **2)** the artifacts brought by the speech enhancement front-end can lead to ASR performance degradation. To this end, a tight integration of the audio-visual speech separation, dereverberation and recognition components via joint fine-tuning [19], [23], [67], [72], [78]–[82] is considered in this paper. Three fine-tuning methods are investigated: **a)** only fine-tuning the back-end ASR component using the enhanced speech outputs while

the front-end remains unchanged; **b)** end-to-end jointly fine-tuning the entire system including the speech enhancement front-end and the recognition back-end components using the ASR cost function; **c)** end-to-end jointly fine-tuning the entire system using a multi-task criterion interpolation between the speech enhancement and recognition cost functions as follows:

$$\mathcal{L} = (1 - \gamma)\mathcal{L}_{\text{ASR}} + \gamma\mathcal{L}_{\text{SE}}, \quad (34)$$

where $\gamma$ is empirically set as 0.5 in the experiments unless otherwise stated. The precise form of the speech enhancement loss function, $\mathcal{L}_{\text{SE}}$, is determined by the underlying integrated front-end architectures being used, as described in Section IV. This is expressed as follows: **a)** $\mathcal{L}_{\text{SE}} = \mathcal{L}_{\text{SISNR}}$ for audio-visual speech separation followed by dereverberation, as in Section IV-A; **b)** $\mathcal{L}_{\text{SE}} = \mathcal{L}_{\text{MSE}}$ for audio-visual speech dereverberation followed by separation, as in Section IV-B; and **c)** $\mathcal{L}_{\text{SE}} = \mathcal{L}_{\text{SISNR}} + \mathcal{L}_{\text{MSE}}$ for joint speech separation & dereverberation in Section IV-C.

## VI. EXPERIMENTAL SETUP

This section is organized as follows. Section VI-A gives the details of the LRS2 corpus. The simulated and replayed multi-channel mixture speech datasets are described in Section Section VI-B and VI-C, respectively. Section VI-D presents the performance of the baseline single-channel ASR and AVSR systems on mixture speech. Finally, two important implementation issues that affect the performance of the proposed audio-visual multi-channel speech separation, dereverberation and recognition systems are discussed in Section VI-E.

### A. LRS2 Corpus

The Oxford LRS2 corpus [84] is one of the largest publicly available corpora for audio-visual speech recognition. This corpus consists of news and talk shows from BBC programs. This is a challenging AVSR task since it contains thousands of speakers with large variations in head pose. The LRS2 corpus is divided into four subsets, i.e. Pre-train, Train, Validation and Test sets. In our experiments, the official Pre-train and Train data sets are combined for model training.

### B. Simulated Overlapped and Reverberant Speech

Since there is no publicly available audio-visual multi-channel mixture speech corpus, we simulated the multi-channel mixture speech with overlapping and reverberation based on the LRS2 corpus in the experiments. Details of the simulation process are described in Algorithm 1. A 15-channel symmetric linear array with non-even inter-channel spacing [7,6,5,4,3,2,1,1,2,3,4,5,6,7]cm is used in the simulation process. 843 point-source noises [99] and 20000 room impulse responses (RIRs) generated by the image method [100] in 400 different simulated rooms are used in our experiment. The distance between a sound source and the microphone array center is uniformly sampled from a range of 1m to 5m and the room size ranges from 4m×4m×3m to 10m×10m ×6m (length×width×height). The reverberation time $T_{60}$ is uniformly sampled from a range of 0.14s to 0.92s. The average overlapping ratio is around 80%. The signal-to-noise ratio (SNR) is uniformly sampled from {0, 5, 10, 15, 20}dB, and the signal-to-interference ratio (SIR) is uniformly sampled

---

**Algorithm 1:** Multi-channel mixture speech simulation

**Input:** single-channel anechoic LRS2 corpus
**Output:** multi-channel mixture speech
**foreach** *utterance* in LRS2 **do**

  1) Uniformly sample an interfering utterance from another speaker in the LRS2 corpus;
  2) Uniformly sample a room size from 4m×4m×3m to 10m×10m×6m;
  3) Uniformly sample a $T_{60}$ from 0.14s to 0.92s;
  4) Uniformly sample a microphone array position in the room;
  5) Uniformly sample two speakers' positions while the distance between each speaker and the array is within the range of 1m to 5m;
  6) Uniformly sample an angle difference from {[0°, 15°), [15°, 45°), [45°, 90°), [90°, 180°) };
  **while** the angle difference of the target and interfering speakers relative to the microphone array not in the selected range **do**
  ⌊ 7) Re-sample the interfering speaker's position;
  8) Generate two multi-channel RIRs for the target and interfering speakers using the above settings and applying the image method [100];
  9) Convolve each single-channel anechoic speech of current utterance with the corresponding multi-channel RIRs to simulate room reverberation;
  10) Uniformly sample a SIR from {-6, 0, 6} dB;
  11) Scale the target and interfering sources with the sampled SIR;
  12) Uniformly sample a noise from a total of 843 point-source noise types [99];
  13) Add two scaled speaker speech signals along with the selected noise under {0, 5, 10, 15, 20}dB SNR to obtain the final multi-channel mixture (overlapped, noisy and reverberant) speech.

---

from {-6, 0, 6}dB. In addition, the angle difference relative to the microphone array between the target and interfering speakers is uniformly sampled from four ranges of the angle difference {[0°, 15°), [15°, 45°), [45°, 90°), [90°, 180°)}. The final simulated multi-channel datasets contain three subsets with 96997, 4272 and 4972 utterances respectively for training (91.37 hours), validation (2.59 hours) and test (2.32 hours).

### C. Replayed Mixture Speech

To further evaluate the performance of the proposed approach in a more realistic application environment, a replayed test set [67] with 1200 utterances (0.5 hours) of LRS2 Test set recorded in a 10m×5m×3m meeting room is also used in our experiments. Two loudspeakers are used to replay different utterances simultaneously to produce mixture speech. The geometric specification of the microphone array used during recording is the same as that used in the simulation. The target and interfering speakers are located at the following directions relative to the microphone array, i.e. {15°/30°, 45°/30°, 75°/30°, 105°/30°, 30°/60°, 90°/60°, 120°/60°, 150°/60°},

TABLE I
PERFORMANCE OF SINGLE-CHANNEL ASR AND AVSR SYSTEMS (WITHOUT SPEECH ENHANCEMENT FRONT-END) TRAINED AND EVALUATED ON ANECHOIC, REVERBERANT-ONLY AND MIXTURE SPEECH. "SIMU" AND "REPLAY" DENOTE THE SIMULATED AND REPLAYED EVALUATION DATASETS OF SECTION VI-B AND SECTION VI-C.

| Sys. | Data | +Visual Features | WER(%) Simu | WER(%) Replay |
|---|---|---|---|---|
| 1 | Anechoic | ✗ | 8.8 | - |
| 2 | Anechoic | ✓ | 7.3 | - |
| 3 | Reverberant-only | ✗ | 13.8 | - |
| 4 | Reverberant-only | ✓ | 10.5 | - |
| 5 | Mixture of raw channel 1 | ✗ | 57.5 | 58.6 |
| 6 | Mixture of raw channel 1 | ✓ | 25.2 | 22.6 |

where the distance between the loudspeakers and microphones ranges from 1m to 1.5m. In the replayed data, the target speaker's DOA is captured by a $180°$ camera [67]. The average overlapping ratio of the replayed mixture speech is around 80% and SIR is around 1.5dB.

### D. Baseline System Description

**1) Speech Enhancement Front-end:** The 257-dimensional complex spectrum of each channel is extracted using a 512-point STFT with a 32ms square-root Hanning window and 16ms frame rate (e.g. Fig. 1, top left corner). The AF and IPD features are computed using 9 microphone pairs {1/15, 2/14, 3/13, 1/7, 12/4, 11/5, 12/8, 7/10, 8/9} to sample different spacing between microphones following [67]. For each Dilated 1D Conv Block in a TCN module (Fig. 2, left), the number of channels in the 1×1 Conv layer is set to 256. The kernel size of the D-Conv layer is set to 3, with 512 channels. The output dimension of the linear layer is set to 257.

**2) Visual Front-end:** The original $160×160$ dimensional video frames in the LRS2 datasets are centrally cropped by a $112×112$ dimensional window and then up-sampled to be time synchronised with the audio frames via linear interpolation. The Visual Front-end (e.g. Fig. 1, bottom left corner, in light green) uses the same hyper-parameter settings as described in [93]. In addition, the number of the acoustic subspaces $K$ is set to 10 with $\mathbf{P}_k^a \in \mathbb{R}^{256×256}$ and $\mathbf{P}^v \in \mathbb{R}^{10×256}$ in the factorized attention layer [36].

**3) Recognition Back-end:** The 80-dimensional log Mel-FBK features extracted using a 25ms window and 10ms frame rate serve as the inputs to the recognition back-end. The baseline Conformer models consist of 12 encoder and 6 decoder blocks following the ESPnet recipe[6]. Each encoder or decoder block is configured with 4-head attention of 256 dimensions and 2048 feed-forward hidden units. The convolutional sub-sampling module includes two 2D convolutional layers with a stride of 2, each followed by a ReLU activation. 500 byte-pair-encoding (BPE) tokens are used as decoder outputs. All models are trained using NVIDIA A40 GPU cards[7].

**4) Performance of Speech Recognition without Speech Enhancement Front-end:** Table I presents the WER results of the single-channel input based Conformer ASR and AVSR

[6]github.com/espnet/espnet/blob/master/egs/lrs2/asr1/run.sh

[7]The jointly fine-tuned speech enhancement front-end and recognition back-end systems in Table V are trained using one thread on a single Nvidia A40 GPU with a batch size of 24 and the GPU memory usage vary from 32G to 43G maximum.

TABLE II
PERFORMANCE OF THREE INTEGRATED SPEECH ENHANCEMENT FRONT-END ARCHITECTURES WITH DIFFERENT NUMBERS OF FILTER TAPS ($L$) ON SIMULATED MIXTURE SPEECH FOR SINGLE-CHANNEL DNN-WPE, MULTI-CHANNEL DNN-WPE AND MASK-BASED WPD MODULES USED IN AUDIO-ONLY SPEECH ENHANCEMENT FRONT-ENDS.

| Sys. | Filter taps ($L$) | PESQ(↑) / STOI(↑) / SRMR(↑) Sep. → Dervb. (Single-channel DNN-WPE) | Dervb. → Sep. (Multi-channel DNN-WPE) | Joint Sep. & Dervb. (Mask-based WPD) |
|---|---|---|---|---|
| 1 | 1 | 2.21/72.07/5.32 | 2.44/79.63/6.31 | **2.42/76.63/6.64** |
| 2 | 2 | 2.22/72.42/5.29 | **2.46/79.75/6.44** | 2.40/76.64/6.83 |
| 3 | 3 | 2.23/72.69/5.32 | 2.45/79.66/6.50 | 2.40/76.51/6.97 |
| 4 | 4 | 2.23/72.86/5.35 | 2.45/79.53/6.57 | 2.36/76.10/7.04 |
| 5 | 5 | 2.24/72.98/5.39 | 2.44/79.32/6.60 | 2.34/75.78/7.08 |
| 6 | 7 | 2.24/73.20/5.45 | 2.41/78.47/6.72 | 2.30/75.05/7.11 |
| 7 | 9 | 2.24/73.35/5.51 | 2.38/77.87/6.70 | 2.27/74.48/7.16 |
| 8 | 12 | 2.25/73.53/5.58 | 2.34/76.73/6.80 | 2.20/73.28/7.12 |
| 9 | 15 | 2.25/73.65/5.64 | 2.28/75.20/6.83 | 2.12/71.74/6.90 |
| 10 | 18 | **2.25/73.73/5.70** | 2.24/74.18/6.90 | 2.06/70.39/6.66 |
| 11 | 21 | 2.25/73.71/5.75 | 2.18/72.67/6.84 | 1.98/68.90/6.48 |
| 12 | 24 | 2.25/73.71/5.79 | 2.11/71.09/6.90 | 1.87/66.20/6.02 |
| 13 | 27 | 2.25/73.70/5.82 | 2.02/68.96/6.72 | 1.81/64.60/5.83 |

systems (without using a microphone array and any speech enhancement front-end) on the anechoic, reverberant-only and mixture speech. It can be observed that using visual information can consistently improve the recognition performance over the audio-only ASR systems by up to **1.5% absolute** (**17.0% relative**) WER reduction on the anechoic speech (sys. 2 vs. sys. 1) and **3.3% absolute** (**23.9% relative**) WER reduction on the reverberant-only speech (sys. 4 vs. sys. 3). In particular, the AVSR system significantly outperforms the audio-only ASR system (sys. 6 vs. sys. 5) by up to **32.3%** and **36.0% absolute** (**56.2%** and **61.4% relative**) WER reductions on the simulated and replayed mixture speech respectively.

### E. Implementation Details

**1) Number of Filter Taps:** The number of filter taps $L$ used in WPE and WPD approaches has a huge impact on the quality of the enhanced speech and the downstream recognition performance. A set of ablation studies on the settings of filter taps $L$ are conducted for each of the three integrated speech separation and dereverberation front-end architectures of Section IV (i.e. "Sep. → Dervb", "Dervb. → Sep." and "Joint Sep. & Dervb." denote the speech separation followed by dereverberation, speech dereverberation followed by separation and joint speech separation & dereverberation, respectively.) These are shown in Table II for audio-only speech enhancement. Considering the speech enhancement performance in terms of PESQ, STOI and SRMR scores, the number of filter taps for single-channel DNN-WPE, multi-channel DNN-WPE and mask-based WPD are respectively chosen and fixed as 18 (sys. 10), 2 (sys. 2) and 1 (sys. 1) in the following experiments. In addition, the prediction delay $D$ is empirically set to 2 for DNN-WPE and mask-based WPD.

**2) Matrix Inversion:** The inversion of the PSD matrices for MVDR and WPD (Eqn. (6) and Eqn. (30)) and the temporal correlation matrix for WPE (Eqn. (15)) are prone to numerical issues when they are ill-conditioned or singular. To this end, the diagonal variance flooring approach [72] is utilized in this work. A complex PSD or correlation matrix $\mathbf{\Phi}$ is floored as $\mathbf{\Phi}' = \mathbf{\Phi} + \varepsilon \operatorname{tr}(\mathbf{\Phi})\mathbf{I}$ before inversion, where a flooring scaling term $\varepsilon$ needs to be set, and $\mathbf{I}$ is the identity matrix. In addition, a more stable complex matrix inversion algorithm [101] is adopted in this paper. A set of ablation studies on the setting of the flooring scaling $\varepsilon$ is shown in Table III for audio-only speech enhancement front-end systems with different

TABLE III
PERFORMANCE OF SPEECH ENHANCEMENT FRONT-ENDS WITH DIFFERENT
DIAGONAL VARIANCE FLOORING ($\varepsilon$) ON SIMULATED MIXTURE SPEECH
FOR MASK-BASED MVDR, SINGLE-CHANNEL DNN-WPE,
MULTI-CHANNEL DNN-WPE AND MASK-BASED WPD USED IN
AUDIO-ONLY SPEECH ENHANCEMENT FRONT-ENDS.

| Sys. | Variance flooring ($\varepsilon$) | PESQ($\uparrow$) / STOI($\uparrow$) / SRMR($\uparrow$) | | | |
|---|---|---|---|---|---|
| | | Sep. (Mask-based MVDR) | Sep. $\rightarrow$ Dervb. (Single-channel DNN-WPE) | Dervb. $\rightarrow$ Sep. (Multi-channel DNN-WPE) | Joint Sep. & Dervb. (Mask-based WPD) |
| 1 | $10^{-1}$ | 1.89/63.41/4.36 | 2.21/71.98/5.67 | 2.36/77.68/6.01 | 2.11/67.85/5.47 |
| 2 | $10^{-3}$ | 2.08/68.50/4.77 | 2.24/73.39/5.88 | 2.44/79.37/6.24 | 2.36/75.38/6.45 |
| 3 | $10^{-4}$ | 2.17/70.39/5.34 | 2.25/73.64/5.80 | 2.43/79.25/6.18 | **2.42/76.63/6.64** |
| 4 | $10^{-5}$ | **2.21/71.30/5.45** | **2.25/73.73/5.70** | 2.45/79.68/6.40 | 1.96/61.45/6.29 |
| 5 | $10^{-6}$ | 2.19/71.24/5.46 | 2.25/73.74/5.65 | **2.46/79.75/6.44** | 1.55/45.13/4.84 |
| 6 | $10^{-7}$ | 2.16/70.92/5.39 | 2.25/73.75/5.64 | 2.44/79.62/6.56 | 1.63/48.27/4.74 |
| 7 | $10^{-9}$ | 1.99/67.02/5.23 | 2.25/73.74/5.64 | 2.25/73.67/5.95 | 1.51/43.99/4.27 |

separation only or integrated (separation and dereverberation) architectures. Based on the PESQ, STOI and SRMR scores, $10^{-5}$ (sys. 4), $10^{-5}$ (sys. 4), $10^{-6}$ (sys. 5) and $10^{-4}$ (sys. 3) are selected as the optimal values of the diagonal variance flooring scaling $\varepsilon$ for mask-based MVDR, single-channel DNN-WPE, multi-channel DNN-WPE and mask-based WPD respectively in the following experiments.

## VII. EXPERIMENTAL RESULTS

In this section, the performance of three integrated audio-visual multi-channel speech separation, dereverberation and recognition architectures of Section IV are evaluated on the LRS2 simulated and replayed mixture speech datasets. Section VII-A analyses the performance improvements by incorporating visual features into different speech enhancement front-end components as well as the recognition back-end. After end-to-end joint fine-tuning, the performance of tightly integrated audio-visual speech separation, dereverberation and recognition systems are presented in Section VII-B.

### A. Performance of Audio-visual Multi-channel Speech Enhancement and Recognition Systems

In this part, we systematically investigate the performance improvements attributed to the visual modality in the proposed integrated speech enhancement architectures of Section IV on the LRS2 simulated multi-channel mixture dataset with four angle difference ranges $[0°, 15°)$, $[15°, 45°)$, $[45°, 90°)$ and $[90°, 180°)$. The mask-based MVDR approach is used in the separation module, and the dereverberation module leverages either DNN-WPE or SpecM based dereverberation methods. The mask-based WPD is used for joint speech separation & dereverberation. The multi-channel audio (including AF and IPD) features and visual modality features and their fusion mechanism presented in Sections II-C, II-D, II-E and III-D for speech separation and dereverberation are used. The visual features are also incorporated into the Conformer speech recognition back-end, as described in Section V. The speech recognition systems in Table IV are obtained by fine-tuning the baseline single-channel Conformer ASR (Table I, sys. 1) or AVSR (Table I, sys. 2) systems using the enhanced outputs of the corresponding speech enhancement front-ends.

From Table IV, several trends can be observed:

**1)** The proposed audio-visual multi-channel speech separation, dereverberation and recognition systems (sys. 11,18,25,32,36) consistently outperformed the corresponding audio-only baseline systems (sys. 5,12,19,26,33) on the LRS2 simulated test set. Consistent performance improvements in PESQ, STOI and SRMR scores were also obtained. For example, a statistically significant WER reduction of **12.4% absolute** (**45.1% relative**) was obtained by the full audio-visual system (sys. 25) over the corresponding audio-only baseline (sys. 19) using a pipelined front-end architecture whereby speech dereverberation was followed by separation. A general trend can also be found that the performance gap between systems with full incorporation of video modality (sys. 11,18,25,32,36) and those using audio-only (sys. 5,12,19,26,33) was much larger when examining the performance on the more challenging subsets, e.g. when inter-speaker angle difference fell in the smallest range of $[0°, 15°)$.

**2)** When compared with audio-only dereverberation, incorporating visual information into the corresponding DNN-WPE (sys. 6,8,10,20,22,24 vs. sys. 5,7,9,19,21,23) or SpecM based dereverberation (sys. 13,15,17,27,29,31 vs. sys. 12,14,16,26,28,30) module produced consistent improvements in terms of PESQ, STOI and SRMR scores, irrespective of the underlying form of integration between speech separation and dereverberation components. A statistically significant WER reduction by up to **1.9% absolute** (sys. 13 vs. sys. 12, **5.9% relative**) was also obtained.

**3)** Among the proposed architectures to integrate speech separation and dereverberation components within the speech enhancement front-end, a pipelined, full audio-visual configuration performing DNN-WPE based speech dereverberation followed by mask-based MVDR speech separation using visual input in both enhancement and recognition stages (sys. 25 vs. sys. 11,18,32,36) produced the lowest overall WERs.

**4)** The integrated audio-visual speech separation, dereverberation and recognition systems (sys. 11,18,25,32,36) consistently outperformed the corresponding separation-only AVSR systems (sys. 4) in terms of PESQ, STOI and SRMR scores. However, with regard to recognition performance, the SpecM based AVSR systems (sys. 18,32) and the mask-WPD based AVSR system (sys. 36) did not outperform the baseline system (sys. 4). The potential causes were: **a)** For systems using SpecM based dereverberation (sys. 18,32), although perceptually enhanced speech quality was obtained when compared to the corresponding baseline systems (sys. 4), the spectral artifacts caused by SpecM introduced a negative impact on downstream speech recognition performance; and **b)** For mask-based WPD systems, the number of filter taps and microphone channels together produced spatial-temporal PSD matrices in Eqns. (31)-(32) larger than, for example, those in Eqns. (7)-(8) for MVDR speech separation only, and thus increased difficulty in their inversion. This was further suggested by the larger variance flooring scaling $\varepsilon=10^{-4}$ in mask-based WPD than all the other systems shown in the ablation studies of Table III. This issue can offset the benefit of joint speech separation & dereverberation from WPD.

**5)** Finally, incorporating both the video modality and AF spatial features into the front-ends (e.g. sys. 3,10,17,24,31,35) consistently outperformed the comparable systems using either only AF features (sys. 1,5,12,19,26,33), or video features alone (sys. 2,8,15,22,29,34).

TABLE IV

PERFORMANCE OF INTEGRATED ARCHITECTURES FOR AUDIO-VISUAL MULTI-CHANNEL SPEECH SEPARATION ("SEP."), DEREVERBERATION ("DERVB.") AND RECOGNITION ("RECG.") ON THE LRS2 SIMULATED MULTI-CHANNEL MIXTURE DATASET. "ARCH.", "AF", "SPECM", "CONF." AND "AVG." DENOTE THE ARCHITECTURE, ANGLE FEATURE, SPECTRAL MAPPING, CONFORMER AND AVERAGE, RESPECTIVELY. $[a°, b°]$ DENOTES THE RANGE OF INTER-SPEAKER ANGLE DIFFERENCE BETWEEN THE TARGET AND INTERFERING SPEAKERS RELATIVE TO THE MICROPHONE ARRAY. "*" AND "†" REPRESENT A STATISTICALLY SIGNIFICANT WER DIFFERENCE OVER THE CORRESPONDING AUDIO-ONLY BASELINE SYSTEMS (SYS. 5,12,19,26,33) AND AUDIO-ONLY DEREVERBERATION BASELINE SYSTEMS (SYS. (5,7,9),(12,14,16),(19,21,23),(26,28,30)), RESPECTIVELY.

| Arch. | Sys. | +AF | Sep. (MVDR) | Dervb. (DNN-WPE) | (SpecM) | Recg. (Conf.) | PESQ(↑) / STOI(↑) / SRMR(↑) / WER(↓) [0°, 15°] | [15°, 45°] | [45°, 90°] | [90°, 180°] | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Mixture of raw channel 1 | | | | | | ✗ | 1.54/53.98/3.58/57.9 | 1.53/53.48/3.57/57.4 | 1.53/53.58/3.58/57.8 | 1.54/54.08/3.60/57.0 | 1.54/53.78/3.58/57.5 |
| | | | | | | ✓ | 1.54/53.98/3.58/25.9 | 1.53/53.48/3.57/24.7 | 1.53/53.58/3.58/25.6 | 1.54/54.08/3.60/24.5 | 1.54/53.78/3.58/25.2 |
| **Sep.** (MVDR only) | 1 | ✓ | ✗ | - | - | ✗ | 1.87/62.07/5.03/51.4 | 2.23/71.90/5.50/29.3 | 2.35/74.95/5.58/22.8 | 2.39/76.28/5.67/21.6 | 2.21/71.30/5.45/31.3 |
| | 2 | ✗ | ✓ | - | - | ✗ | 2.20/71.35/5.34/28.6 | 2.29/73.39/5.52/24.9 | 2.33/74.41/5.52/23.9 | 2.35/75.19/5.60/22.4 | 2.29/73.59/5.50/25.0 |
| | 3 | ✓ | ✓ | - | - | ✗ | 2.15/70.30/5.32/33.1 | 2.30/73.95/5.63/23.0 | 2.38/75.71/5.65/21.4 | 2.42/76.96/5.74/19.7 | 2.31/74.23/5.59/24.3 |
| | 4 | ✓ | ✓ | - | - | ✓ | 2.15/70.30/5.32/21.7 | 2.30/73.95/5.63/15.9 | 2.38/75.71/5.65/14.7 | 2.42/76.96/5.74/13.2 | 2.31/74.23/5.59/16.4 |
| **Sep. → Dervb.** (MVDR → DNN-WPE) | 5 | ✓ | ✗ | ✗ | | ✗ | 1.91/64.21/5.23/50.1 | 2.26/74.36/5.74/27.9 | 2.39/77.51/5.86/21.6 | 2.44/78.82/5.99/20.5 | 2.25/73.73/5.70/30.0 |
| | 6 | ✓ | ✗ | ✓ | | ✗ | 1.91/64.51/5.22/49.6 | 2.26/74.56/5.76/27.4 | 2.39/77.69/5.87/21.1 | 2.44/78.96/6.01/20.3 | 2.25/73.93/5.71/29.6† |
| | 7 | ✗ | ✓ | ✗ | | ✗ | **2.24**/73.72/5.57/28.5 | 2.33/75.84/5.79/25.2 | 2.37/76.95/5.80/22.9 | 2.40/77.65/5.94/22.2 | 2.34/76.04/5.78/24.7 |
| | 8 | ✗ | ✓ | ✓ | | ✗ | **2.24/73.81/5.60**/28.2 | 2.33/75.90/5.80/24.5 | 2.38/77.03/5.82/23.2 | 2.40/77.74/5.95/22.3 | 2.34/76.12/5.79/24.5 |
| | 9 | ✓ | ✓ | ✗ | | ✗ | 2.18/72.69/5.54/31.9 | **2.35**/76.53/**5.92**/23.1 | **2.43**/78.43/5.94/20.2 | **2.47**/79.63/**6.09**/18.9 | **2.36**/76.82/5.87/23.5 |
| | 10 | ✓ | ✓ | ✓ | | ✗ | 2.18/72.79/5.55/31.7 | **2.35/76.62/5.92**/22.7 | 2.42/**78.50/5.95**/20.0 | **2.47/79.67/6.09**/18.2 | **2.36/76.90/5.88**/23.2† |
| | 11 | ✓ | ✓ | ✓ | | ✓ | 2.18/72.79/5.55/**21.1** | **2.35/76.62/5.92/15.2** | 2.42/**78.50**/5.95/**14.1** | **2.47/79.67/6.09**/**13.5** | **2.36/76.90/5.88/16.0*** |
| **Sep. → Dervb.** (MVDR → SpecM) | 12 | ✓ | ✗ | | ✗ | ✗ | 1.95/66.13/7.00/52.7 | 2.37/77.53/7.44/30.7 | 2.51/80.81/7.50/23.6 | 2.57/81.91/7.54/22.5 | 2.35/76.60/7.37/32.4 |
| | 13 | ✓ | ✗ | | ✓ | ✗ | 1.98/67.73/7.16/50.6 | 2.41/78.67/7.66/28.3 | 2.55/81.81/7.71/22.8 | 2.60/82.77/7.73/20.4 | 2.39/77.75/7.56/30.5† |
| | 14 | ✗ | ✓ | | ✗ | ✗ | 2.37/78.10/7.50/31.6 | 2.47/80.10/7.65/27.3 | 2.52/81.12/7.63/24.9 | 2.54/81.60/7.66/24.1 | 2.48/80.23/7.61/27.0 |
| | 15 | ✗ | ✓ | | ✓ | ✗ | **2.38/78.36/7.62**/29.9 | 2.48/80.28/7.71/25.3 | 2.53/81.25/7.71/23.8 | 2.55/81.69/7.73/22.7 | 2.49/80.39/7.69/25.4† |
| | 16 | ✓ | ✓ | | ✗ | ✗ | 2.31/76.74/7.43/35.0 | 2.48/80.72/7.72/25.8 | **2.60**/82.61/7.68/22.3 | 2.64/**83.55**/7.73/20.9 | **2.51**/80.91/7.74/26.0 |
| | 17 | ✓ | ✓ | | ✓ | ✗ | 2.31/76.98/7.58/33.5 | **2.51/80.86/7.84**/23.8 | 2.59/**82.71/7.80**/21.7 | 2.64/83.55/**7.85**/19.2 | **2.51/81.03/7.77**/24.5† |
| | 18 | ✓ | ✓ | | ✓ | ✓ | 2.31/76.98/7.58/**22.0** | **2.51/80.86/7.84/16.8** | 2.59/**82.71/7.80**/**14.5** | **2.64/83.55/7.85**/14.4 | **2.51/81.03/7.77/16.9*** |
| **Dervb. → Sep.** (DNN-WPE → MVDR) | 19 | ✓ | ✗ | ✗ | | ✗ | 2.04/69.13/5.86/47.2 | 2.48/80.58/6.46/24.6 | 2.63/84.17/6.68/19.2 | 2.68/85.11/6.75/19.2 | 2.46/79.75/6.44/27.5 |
| | 20 | ✓ | ✗ | ✓ | | ✗ | 2.03/69.57/5.85/46.7 | 2.46/80.46/6.39/24.6 | 2.62/83.93/6.66/19.2 | 2.67/85.14/6.78/19.1 | 2.45/79.78/6.42/27.4 |
| | 21 | ✗ | ✓ | ✗ | | ✗ | 2.39/78.99/6.26/27.4 | 2.53/81.61/6.61/22.0 | 2.60/83.23/6.69/20.4 | 2.61/83.15/6.71/20.9 | 2.53/81.75/6.57/22.7 |
| | 22 | ✗ | ✓ | ✓ | | ✗ | **2.41/79.37/6.33**/25.8 | 2.54/81.90/6.64/21.4 | 2.61/83.33/6.67/19.6 | 2.61/83.40/6.74/20.4 | 2.54/82.00/6.59/21.8† |
| | 23 | ✓ | ✓ | ✗ | | ✗ | 2.34/77.63/6.19/30.3 | **2.57/82.62/6.65**/20.2 | **2.68**/84.89/**6.81**/17.9 | **2.71**/85.64/**6.89**/17.1 | **2.57/82.69/6.64**/21.4 |
| | 24 | ✓ | ✓ | ✓ | | ✗ | 2.33/77.39/6.20/30.7 | 2.55/82.40/6.63/20.3 | 2.66/**84.93**/6.74/17.7 | 2.70/**85.71**/6.85/17.2 | 2.56/82.61/6.61/21.4 |
| | 25 | ✓ | ✓ | ✓ | | ✓ | 2.33/77.39/6.20/**20.8** | 2.55/82.40/6.63/**15.1** | 2.66/**84.93**/6.74/**12.4** | 2.70/**85.71**/6.85/**12.3** | 2.56/82.61/6.61/**15.1*** |
| **Dervb. → Sep.** (SpecM → MVDR) | 26 | ✓ | ✗ | | ✗ | ✗ | 1.82/63.34/5.96/57.0 | 2.22/73.58/6.44/32.8 | 2.43/78.33/6.83/25.6 | 2.49/79.63/6.94/24.2 | 2.24/73.72/6.54/34.9 |
| | 27 | ✓ | ✗ | | ✓ | ✗ | 1.82/62.99/5.82/57.7 | 2.24/74.06/6.43/31.9 | 2.43/78.55/6.82/24.8 | 2.49/79.73/6.96/23.8 | 2.24/73.83/6.51/34.6 |
| | 28 | ✗ | ✓ | | ✗ | ✗ | **2.17**/72.66/6.39/38.8 | 2.33/76.11/6.73/29.3 | 2.45/78.55/6.89/24.9 | 2.47/78.51/6.92/26.4 | 2.35/76.46/6.73/29.8 |
| | 29 | ✗ | ✓ | | ✓ | ✗ | 2.16/**72.82/6.44**/36.0 | 2.31/75.88/6.72/29.3 | 2.43/78.39/6.90/25.1 | 2.41/78.16/6.90/25.1 | 2.33/76.31/6.74/28.9† |
| | 30 | ✓ | ✓ | | ✗ | ✗ | 2.12/71.49/6.28/39.7 | **2.35**/76.58/**6.74**/27.1 | **2.49/79.62/6.94**/22.4 | **2.53/80.73**/7.03/22.8 | **2.37/77.11**/6.75/28.0 |
| | 31 | ✓ | ✓ | | ✓ | ✗ | 2.11/71.59/6.35/38.2 | 2.34/**76.65**/6.71/27.6 | 2.47/79.54/6.92/22.4 | 2.52/80.61/**7.04**/21.1 | 2.36/77.10/**6.76**/27.3† |
| | 32 | ✓ | ✓ | | ✓ | ✓ | 2.11/71.59/6.35/**24.9** | 2.34/**76.65**/6.71/**16.9** | 2.47/79.54/6.92/**15.3** | 2.52/80.61/**7.04**/14.8 | 2.36/77.10/**6.76/18.0*** |
| **Joint Sep. & Dervb.** (WPD) | 33 | ✓ | | ✗ | | ✗ | 1.99/65.92/6.06/55.0 | 2.41/76.81/6.60/30.1 | 2.60/81.09/6.83/22.6 | 2.67/82.69/7.07/21.8 | 2.42/76.63/6.64/32.4 |
| | 34 | ✗ | | ✓ | | ✗ | **2.29/74.79/6.44**/34.3 | 2.46/78.21/**6.78**/25.5 | 2.57/80.26/6.91/22.9 | 2.57/80.48/7.01/22.9 | 2.47/78.44/6.78/26.4 |
| | 35 | ✓ | | ✓ | | ✗ | 2.26/73.78/6.36/37.7 | **2.50/79.11**/6.75/25.1 | **2.64/82.25/6.95**/20.7 | **2.70/83.43/7.13**/20.1 | **2.53/79.64/6.80**/25.9 |
| | 36 | ✓ | | ✓ | | ✓ | 2.26/73.78/6.36/**24.6** | **2.50/79.11**/6.75/**16.3** | **2.64/82.25/6.95**/**13.7** | **2.70/83.43/7.13**/**13.6** | **2.53/79.64/6.80/17.0*** |

## B. Performance of End-to-end Joint Fine-tuning of Speech Enhancement Front-end and Recognition Back-end

The most representative subset of audio-visual and audio-only multi-channel systems in Table IV are then end-to-end joint fine-tuning using either the ASR cost function alone, or a multi-task criterion interpolation between the speech enhancement and recognition cost as described in Section V-B. Their performance in terms of WER and front-end metrics (PESQ, STOI and SRMR) are evaluated on both the LSR2 simulated ("Simu") and replayed ("Replay") test sets and shown in Table V (original system numbering in Table IV carried over). Several main trends can be observed:

**1)** After end-to-end joint fine-tuning, consistent performance improvements in WER were obtained over all systems without doing so (sys. marked with "-" in Col. 3, Table V), irrespective of the joint fine-tuning criterion based on ASR loss alone (sys. marked with "(a)"), or its interpolation with enhancement loss (sys. marked with "(b)"). In particular, statistically significant overall ("O.V.") WER reductions of **3.3%** and **1.6% absolute** (**14.6%** and **11.9% relative**) were obtained using the joint fine-tuned ASR (sys. 19(a) vs. sys. 19) and AVSR (sys. 25(b) vs. sys. 25) systems across both test sets. Consistent performance improvements in speech enhancement front-end metrics

scores were also obtained. Fig. 5 shows a set of example spectra of **(a)** Overlapped-reverberant-noisy speech, **(b)** Target clean speech, **(c)** Pipelined audio-only speech enhancement output (Table IV, sys. 19), **(d)** Pipelined audio-visual speech enhancement output (Table IV, sys. 25), **(e)** Jointly fine-tuned audio-only speech enhancement output (Table V, sys. 19(b)), and **(f)** Jointly fine-tuned audio-visual speech enhancement output (Table V, sys. 25(b)). The spectrum portions circled using blue dotted lines in **(a)** represent the interfering speaker's speech, background noise and reverberation, which have been largely removed in **(f)**.

**2)** The best overall performance was produced by the end-to-end joint fine-tuned audio-visual system with DNN-WPE based dereverberation followed by mask-based MVDR (sys.25(b)). Using this system statistically significant WER reductions of up to **9.1%** and **6.2% absolute** (**41.7%** and **36.0% relative**) were obtained on the LRS2 simulated and replayed test sets over the audio-only baseline (19(b)). In addition, all the jointly fine-tuned audio-visual speech separation, dereverberation and recognition systems consistently outperformed the comparable baseline separation-only AVSR systems (e.g. sys. 11(b),18(b),25(b),32(b),36(b) vs. sys. 4(b)), with a statistically significant WER reduction up to **1.9% absolute** (**13.8% relative**) (sys. 25(b) vs. sys. 4(b)).

TABLE V
PERFORMANCE OF AUDIO-VISUAL AND AUDIO-ONLY MULTI-CHANNEL SPEECH RECOGNITION SYSTEMS AFTER END-TO-END JOINT FINE-TUNING USING ASR COST $\mathcal{L}_{\text{ASR}}$ ALONE (MARKED WITH "(a)"), OR ITS INTERPOLATED WITH ENHANCEMENT LOSS $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ (MARKED WITH "(b)"), ON THE LRS2 SIMULATED ("SIMU") AND REPLAYED ("REPLAY") TEST SETS. THE ORIGINAL SYSTEM NUMBERING FROM TABLE IV IS USED. "AVG." IS IN SHORT FOR "AVERAGE" AND "O.V." FOR "OVERALL" RESULTS ON BOTH SIMULATED AND REPLAYED TEST DATA. "†", "∗" AND "‡" DENOTE A STATISTICALLY SIGNIFICANT WER DIFFERENCE OVER THE SYSTEMS WITHOUT JOINT FINE-TUNING (MARKED WITH "-"), THE CORRESPONDING AUDIO-ONLY BASELINE SYSTEMS (SYS. 5(b), 12(b), 19(b), 26(b), 33(b)) AND SEPARATION-ONLY AVSR BASELINE SYSTEM (SYS. 4(b)), RESPECTIVELY.

| Arch. | Sys. | Jointly Fine-tuning Criterion | PESQ(↑) / STOI(↑) / SRMR(↑) Avg. Simu | Replay | WER(↓) [0°,15°] Simu | Replay | [15°,45°] Simu | Replay | [45°,90°] Simu | Replay | [90°,180°] Simu | Avg. Simu | Replay | O.V. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Sep.** (MVDR only) | 1 | - | 2.21/71.30/5.45 | 2.32/77.77/4.31 | 51.4 | 30.6 | 29.3 | 23.6 | 22.8 | 18.5 | 21.6 | 31.3 | 23.4 | 27.4 |
| | 1(a) | $\mathcal{L}_{\text{ASR}}$ | 2.46/77.72/6.27 | 2.55/81.90/5.35 | 41.5 | 33.0 | 21.1 | 20.7 | 17.0 | 18.2 | 16.2 | 24.0 | 22.8 | 23.4 |
| | 1(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.32/74.75/5.77 | 2.40/80.11/4.61 | 42.2 | 28.7 | 22.8 | 20.1 | 18.2 | 17.8 | 17.9 | 25.3 | 21.4 | 23.4 |
| | 4 | - | 2.31/74.23/5.59 | 2.37/79.18/4.42 | 21.7 | 15.9 | 15.9 | 12.8 | 14.7 | 13.6 | 13.2 | 16.4 | 13.9 | 15.2 |
| | 4(a) | $\mathcal{L}_{\text{ASR}}$ | 2.53/79.68/6.39 | 2.58/83.47/5.51 | 17.0 | 15.5 | 13.2 | 11.8 | 11.7 | 11.0 | 11.4 | 13.3 | 12.4 | 12.9 |
| | 4(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.38/76.36/5.77 | 2.42/80.81/4.60 | 18.5 | 15.8 | 14.4 | 12.4 | 12.6 | 12.1 | 12.2 | 14.4 | 13.1 | 13.8 |
| **Sep. → Dervb.** (MVDR → DNN-WPE) | 5 | - | 2.25/73.73/5.70 | 2.41/80.19/4.86 | 50.1 | 27.5 | 27.9 | 22.3 | 21.6 | 16.4 | 20.5 | 30.0 | 21.4 | 25.7 |
| | 5(a) | $\mathcal{L}_{\text{ASR}}$ | 2.46/78.62/6.46 | 2.58/82.96/5.80 | 39.9 | 27.4 | 20.4 | 18.9 | 16.2 | 17.7 | 15.8 | 23.1† | 20.6 | 21.9† |
| | 5(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.45/78.27/6.42 | 2.56/82.39/5.72 | 40.7 | 31.8 | 20.8 | 20.1 | 16.4 | 17.8 | 16.0 | 23.5† | 22.1 | 22.8† |
| | 11 | - | 2.36/76.90/5.88 | 2.46/81.71/5.04 | 21.1 | 15.7 | 15.2 | 12.5 | 14.1 | 12.1 | 13.5 | 16.0 | 13.2 | 14.6 |
| | 11(a) | $\mathcal{L}_{\text{ASR}}$ | 2.58/81.26/6.61 | 2.67/84.82/6.14 | **16.7** | 12.9 | **12.8** | 12.2 | 11.7 | 10.8 | 11.0 | **13.0†** | 11.9† | 12.5† |
| | 11(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.55/80.69/6.66 | 2.66/84.77/6.15 | 17.2 | **12.0** | 13.2 | **11.6** | 11.8 | 10.3 | 11.2 | 13.3†*‡ | **11.2†*‡** | **12.3†*‡** |
| **Sep. → Dervb.** (MVDR → SpecM) | 12 | - | 2.35/76.60/7.37 | 2.48/80.75/6.62 | 52.7 | 31.2 | 30.7 | 25.3 | 23.6 | 19.6 | 22.5 | 32.4 | 24.6 | 28.5 |
| | 12(a) | $\mathcal{L}_{\text{ASR}}$ | 2.52/79.84/7.23 | 2.61/83.59/6.59 | 38.3 | 30.4 | 20.6 | 19.8 | 16.3 | 16.4 | 15.9 | 22.8† | 21.2† | 22.0† |
| | 12(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.49/79.16/6.61 | 2.56/82.99/5.91 | 39.4 | 30.3 | 20.9 | 19.4 | 16.2 | 17.1 | 16.2 | 23.2† | 21.2† | 22.2† |
| | 18 | - | 2.51/81.03/**7.77** | 2.58/82.99/**7.29** | 22.0 | 17.2 | 16.8 | 13.4 | 14.5 | 13.1 | 14.4 | 16.9 | 14.2 | 15.6 |
| | 18(a) | $\mathcal{L}_{\text{ASR}}$ | **2.60**/81.64/7.41 | **2.68**/**85.22**/6.77 | **16.7** | 13.8 | 13.0 | 12.2 | **11.6** | 10.8 | 11.0 | **13.1†** | 12.1† | 12.6† |
| | 18(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.55/80.65/6.70 | 2.60/84.43/6.05 | **16.7** | 13.4 | 13.0 | 12.3 | 11.9 | 10.6 | **10.9** | 13.1†*‡ | 11.9†*‡ | 12.5†*‡ |
| **Dervb. → Sep.** (DNN-WPE → MVDR) | 19 | - | 2.46/79.75/6.44 | 2.67/84.68/6.32 | 47.2 | 25.4 | 24.6 | 15.6 | 19.2 | 13.2 | 19.2 | 27.5 | 17.1 | 22.6 |
| | 19(a) | $\mathcal{L}_{\text{ASR}}$ | 2.61/81.91/6.86 | 2.70/85.21/6.28 | 37.8 | 22.2 | 18.8 | 17.2 | 14.9 | 13.1 | 15.0 | 21.6† | 16.9 | 19.3† |
| | 19(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.61/82.12/6.82 | 2.69/85.22/6.28 | 37.6 | 25.3 | 19.0 | 15.5 | 15.6 | 13.5 | 15.0 | 21.8† | 17.2 | 19.5† |
| | 25 | - | 2.56/82.61/6.61 | 2.72/85.85/6.49 | 20.8 | 15.0 | 15.1 | 12.0 | 12.4 | 10.7 | 12.3 | 15.1 | 11.8 | 13.5 |
| | 25(a) | $\mathcal{L}_{\text{ASR}}$ | **2.71**/84.33/**7.04** | **2.75**/86.42/6.48 | **16.0** | 14.4 | 12.5 | **10.6** | **10.7** | 10.1 | 11.2 | **12.6†** | 11.4 | **12.0†** |
| | 25(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.68/**84.75**/6.80 | **2.75**/**86.82**/6.50 | 16.2 | 13.3 | 12.7 | **10.6** | 11.0 | 9.9 | 10.8 | 12.7†*‡ | **11.0*‡** | **11.9†*‡** |
| **Dervb. → Sep.** (SpecM → MVDR) | 26 | - | 2.24/73.72/6.54 | 2.51/80.67/6.32 | 57.0 | 30.4 | 32.8 | 20.4 | 25.6 | 14.9 | 24.2 | 34.9 | 20.8 | 27.9 |
| | 26(a) | $\mathcal{L}_{\text{ASR}}$ | 2.52/79.11/6.46 | 2.62/82.60/5.67 | 41.4 | 32.9 | 22.3 | 19.0 | 16.9 | 16.9 | 15.6 | 24.1† | 21.7 | 22.9† |
| | 26(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.53/79.57/6.50 | 2.65/83.12/5.91 | 42.5 | 29.8 | 21.7 | 19.9 | 16.7 | 16.5 | 16.1 | 24.3† | 21.1 | 22.7† |
| | 32 | - | 2.36/77.10/6.76 | 2.57/82.20/**6.60** | 24.9 | 15.5 | 16.9 | 13.2 | 15.3 | 11.2 | 14.8 | 18.0 | 13.0 | 15.5 |
| | 32(a) | $\mathcal{L}_{\text{ASR}}$ | 2.65/82.02/6.76 | 2.68/84.55/5.87 | 16.7 | 14.4 | 12.7 | 11.8 | 11.1 | 10.7 | 10.6 | **12.8†** | 12.1† | 12.5† |
| | 32(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.66/82.94/6.63 | 2.71/85.34/5.98 | 16.8 | 13.4 | **12.4** | 11.5 | 11.0 | 11.4 | **10.3** | 12.6†*‡ | 11.9†*‡ | 12.3†*‡ |
| **Joint Sep. & Dervb.** (WPD) | 33 | - | 2.42/76.63/6.64 | 2.62/83.25/6.12 | 55.0 | 28.8 | 30.1 | 17.4 | 22.6 | 15.0 | 21.8 | 32.4 | 19.4 | 25.9 |
| | 33(a) | $\mathcal{L}_{\text{ASR}}$ | 2.52/78.55/6.97 | 2.63/82.88/6.18 | 43.6 | 34.1 | 23.3 | 14.9 | 17.4 | 17.9 | 17.0 | 25.3† | 20.8 | 23.1† |
| | 33(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.53/78.76/6.95 | 2.64/83.23/6.17 | 44.7 | 32.5 | 23.3 | 15.0 | 18.1 | 17.2 | 17.0 | 25.7† | 20.2 | 23.0† |
| | 36 | - | 2.53/79.64/6.80 | 2.67/84.45/6.29 | 24.6 | 15.8 | 16.3 | 12.1 | 13.7 | 11.3 | 13.6 | 17.0 | 12.7 | 14.9 |
| | 36(a) | $\mathcal{L}_{\text{ASR}}$ | **2.61**/80.93/**7.16** | 2.69/84.81/**6.39** | 19.0 | **12.7** | **13.8** | 11.0 | **11.4** | 10.2 | 11.3 | **13.9†** | **11.1†** | **12.5†** |
| | 36(b) | $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ | 2.60/**81.27**/6.95 | **2.70**/**85.15**/6.34 | 19.4 | 13.7 | 13.9 | **10.5** | 11.9 | 10.7 | 11.5 | 14.2†* | 11.4†*‡ | 12.8†*‡ |



Fig. 5. Example spectra of **(a)** Overlapped-reverberant-noisy speech, **(b)** Target clean speech, **(c)** Pipelined audio-only speech enhancement output (Table IV, sys. 19), **(d)** Pipelined audio-visual speech enhancement output (Table IV, sys. 25), **(e)** Jointly fine-tuned audio-only speech enhancement output (Table V, sys. 19(b)), and **(f)** Jointly fine-tuned audio-visual speech enhancement output (Table V, sys. 25(b)). The spectrum portions circled using blue dotted lines in **(a)** represent the interfering speaker's speech, background noise and reverberation, which have been largely removed in **(f)**.

**3)** End-to-end joint fine-tuning of the speech enhancement front-end and recognition back-end is effective in mitigating the impact from spectral artifacts produced in SpecM based dereverberation [82] (e.g. sys. 12(b),18(b),26(b),32(b)). This leads to their smaller performance gap against systems using DNN-WPE dereverberation (sys. 5(b),11(b),19(b),25(b)), when compared the gap before joint fine-tuning.

**4)** A further ablation study is conducted on the setting of the speech enhancement cost weight $\gamma$ in Eqn. (34) using three end-to-end joint fine-tuned multi-channel speech enhancement and recognition systems: sys. 1(b), 4(b) and 25(b) of Table V. Their WER performance with respect to $\gamma$ on the LRS2 simulated ("Simu") and replayed ("Replay") test sets are shown in Table VI. These results suggest that the performance of the audio-visual multi-channel speech separation, dereverberation and recognition system (sys. 25(b)) is largely insensitive to the setting of $\gamma \in [0, 0.75]$ during end-to-end joint fine-tuning using interpolated speech enhancement and ASR error costs.

TABLE VI
WER(%) PERFORMANCE OF END-TO-END JOINT FINE-TUNED MULTI-CHANNEL SPEECH ENHANCEMENT AND RECOGNITION SYSTEMS 1(b), 4(b) AND 25(b) OF TABLE V WITH RESPECT TO THE SPEECH ENHANCEMENT COST WEIGHT $\gamma$ IN EQN. (34) ON THE LRS2 SIMULATED ("SIMU") AND REPLAYED ("REPLAY") TEST SETS.

| $\gamma$ / Sys. | 0 Simu | Replay | 0.25 Simu | Replay | 0.5 Simu | Replay | 0.75 Simu | Replay | 1 Simu | Replay |
|---|---|---|---|---|---|---|---|---|---|---|
| 1(b) | 24.0 | 22.8 | 24.7 | 22.4 | 25.3 | 21.4 | 27.2 | 22.7 | 31.3 | 23.4 |
| 4(b) | 13.3 | 12.4 | 14.0 | 12.6 | 14.4 | 13.1 | 14.6 | 13.2 | 16.4 | 13.9 |
| 25(b) | **12.6** | **11.4** | 12.7 | 11.3 | 12.7 | 11.0 | 12.7 | 10.4 | 15.1 | 11.8 |

**5)** The performance of the most important systems shown in Table IV (sys. 1,4,5,11,12,18,19,25,26,32,33,36) and Table V (sys. 1(b),4(b),5(b),11(b),12(b),18(b),19(b),25(b),26(b),32(b), 33(b),36(b)) are further evaluated on the LRS3 [102] test set after applying the same multi-channel mixture speech

TABLE VII
PERFORMANCE OF INTEGRATED ARCHITECTURES FOR AUDIO-VISUAL MULTI-CHANNEL SPEECH SEPARATION ("SEP."), DEREVERBERATION ("DERVB.") AND RECOGNITION ("RECG.") ON THE LRS3 TEST SET SIMULATED MULTI-CHANNEL MIXTURE SPEECH VIA THE LRS2 DATA TRAINED PIPELINED AND JOINTLY FINE-TUNED (USING $\mathcal{L}_{\text{ASR}} + \mathcal{L}_{\text{SE}}$ COST FUNCTION) SYSTEMS IN TABLE IV AND TABLE V, RESPECTIVELY. "ARCH.", "AF", "SPECM", "CONF." AND "AVG." DENOTE THE ARCHITECTURE, ANGLE FEATURE, SPECTRAL MAPPING, CONFORMER AND AVERAGE, RESPECTIVELY. "†", "∗" AND "‡" DENOTE A STATISTICALLY SIGNIFICANT WER DIFFERENCE OVER THE "PIPELINED" SYSTEMS, THE CORRESPONDING AUDIO-ONLY BASELINE SYSTEMS (SYS. 5,12,19,26,33) IN THE "JOINTLY FINE-TUNED" COLUMN AND SEPARATION-ONLY AVSR BASELINE SYSTEM (SYS. 4) IN THE "JOINTLY FINE-TUNED" COLUMN, RESPECTIVELY.

| Arch. | Sys. | +AF | +Visual Features Sep. (MVDR) | Dervb. (DNN-WPE) | (SpecM) | Recg. (Conf.) | PESQ(↑) / STOI(↑) / SRMR(↑) / WER(↓) Avg. Pipelined | Jointly fine-tuned |
|---|---|---|---|---|---|---|---|---|
| **Sep.** (MVDR only) | 1 | ✓ | ✗ | - | | ✗ | 2.22/72.63/5.76/40.3 | 2.32/75.77/6.09/34.5 |
|  | 4 | ✓ | ✓ | - | | ✓ | 2.30/75.18/5.89/29.8 | 2.38/77.57/6.12/26.9 |
| **Sep. → Dervb.** (MVDR → DNN-WPE) | 5 | ✓ | ✗ | ✗ | | ✗ | 2.25/74.97/6.12/38.6 | 2.46/79.44/6.95/31.9† |
|  | 11 | ✓ | ✓ | ✓ | | ✓ | 2.34/77.71/6.28/29.5 | 2.55/81.64/7.24/25.1†*‡ |
| **Sep. → Dervb.** (MVDR → SpecM) | 12 | ✓ | ✗ | | ✗ | ✗ | 2.38/77.88/8.02/41.9 | 2.50/80.32/7.13/31.7† |
|  | 18 | ✓ | ✓ | | ✓ | ✓ | 2.51/81.14/8.46/31.1 | 2.56/81.81/7.17/25.3†*‡ |
| **Dervb. → Sep.** (DNN-WPE → MVDR) | 19 | ✓ | ✗ | ✗ | | ✗ | 2.48/81.40/7.25/34.6 | 2.66/83.88/**7.80**/28.9† |
|  | 25 | ✓ | ✓ | ✓ | | ✓ | 2.55/83.22/7.32/27.2 | **2.69/85.73**/7.70/**23.9**†*‡ |
| **Dervb. → Sep.** (SpecM → MVDR) | 26 | ✓ | ✗ | | ✗ | ✗ | 2.28/75.91/7.33/42.4 | 2.54/81.00/7.14/32.5† |
|  | 32 | ✓ | ✓ | | ✓ | ✓ | 2.35/77.73/7.37/32.2 | 2.61/83.16/7.20/25.6†*‡ |
| **Joint Sep. & Dervb.** (WPD) | 33 | ✓ | | ✗ | | ✗ | 2.45/78.84/7.27/39.4 | 2.54/80.46/7.46/34.1† |
|  | 36 | ✓ | | ✓ | | ✓ | 2.51/80.60/7.31/30.3 | 2.58/82.19/7.46/26.7†* |

simulation protocol of Algorithm 1. These results are shown in Table VII. Similar trends of WER reductions and improvements on speech enhancement metric scores, as well as the same performance ranking among the corresponding systems previously shown in Table IV and Table V, can also be found in Table VII.

## VIII. CONCLUSION

In this paper, an audio-visual multi-channel speech separation, dereverberation and recognition approach featuring a full incorporation of visual information into all system components is proposed. The advantages of additional visual modality over using acoustic features only are demonstrated consistently in mask-based MVDR speech separation, DNN-WPE or spectral mapping (SpecM) based speech dereverberation front-end and Conformer based ASR back-end. A set of audio-visual front-end architectures that integrates the speech separation and dereverberation modules in a pipelined or joint fashion are also derived. They are end-to-end jointly fine-tuned to minimize the error cost mismatch between the speech enhancement front-end and ASR back-end. Experiments were conducted on the mixture overlapped and reverberant speech data constructed using simulation or replay of the benchmark Oxford LRS2 dataset. The proposed audio-visual multi-channel speech separation, dereverberation and recognition systems consistently outperformed the comparable audio-only multi-channel baseline by 9.1% and 6.2% absolute (41.7% and 36.0% relative) in word error rate (WER) reductions, together with consistent improvements obtained on PESQ, STOI and SRMR based speech enhancement metrics. Future research will focus on improving system generalization to diverse microphone array geometrics and room acoustics.

## ACKNOWLEDGMENT

## REFERENCES

[1] J. H. McDermott, "The cocktail party problem," *CURR BIOL*, 2009.

[2] Y. Qian *et al.*, "Past review, current progress, and challenges ahead on the cocktail party problem," *FRONT INFORM TECH EL*, 2018.

[3] T. Yoshioka *et al.*, "The NTT CHiME-3 system: Advances in speech enhancement and recognition for mobile multi-microphone devices," in *ASRU*, 2015.

[4] X. Chang *et al.*, "MIMO-Speech: End-to-end multi-channel multi-speaker speech recognition," in *ASRU*, 2019.

[5] R. Haeb-Umbach *et al.*, "Far-field automatic speech recognition," *Proceedings of the IEEE*, 2020.

[6] X. Anguera *et al.*, "Acoustic beamforming for speaker diarization of meetings," *IEEE-ACM T AUDIO SPE*, 2007.

[7] M. Souden *et al.*, "On optimal frequency-domain multichannel linear filtering for noise reduction," *IEEE-ACM T AUDIO SPE*, 2009.

[8] H. L. Van Trees, *Optimum array processing: Part IV of detection, estimation, and modulation theory*. John Wiley & Sons, 2002.

[9] E. Warsitz *et al.*, "Blind acoustic beamforming based on generalized eigenvalue decomposition," *IEEE-ACM T AUDIO SPE*, 2007.

[10] F. Bahmaninezhad *et al.*, "A comprehensive study of speech separation: spectrogram vs waveform separation," in *INTERSPEECH*, 2019.

[11] L. Chen *et al.*, "Multi-band pit and model integration for improved multi-channel speech separation," in *ICASSP*, 2019.

[12] R. Gu *et al.*, "Enhancing end-to-end multi-channel speech separation via spatial feature learning," in *ICASSP*, 2020.

[13] T. N. Sainath *et al.*, "Multichannel signal processing with deep neural networks for automatic speech recognition," *IEEE-ACM T AUDIO SPE*, 2017.

[14] X. Xiao *et al.*, "Deep beamforming networks for multi-channel speech recognition," in *ICASSP*, 2016.

[15] T. Yoshioka *et al.*, "Recognizing overlapped speech in meetings: A multichannel separation approach using neural networks," in *INTERSPEECH*, 2018.

[16] T. Yoshioka *et al.*, "Multi-microphone neural speech separation for far-field multi-talker speech recognition," in *ICASSP*, 2018.

[17] T. Higuchi *et al.*, "Frame-by-frame closed-form update for mask-based adaptive mvdr beamforming," in *ICASSP*, 2018.

[18] Y. Kubo *et al.*, "Mask-based mvdr beamformer for noisy multisource environments: Introduction of time-varying spatial covariance model," in *ICASSP*, 2019.

[19] Y. Xu *et al.*, "Joint training of complex ratio mask based beamformer and acoustic model for noise robust ASR," in *ICASSP*, 2019.

[20] J. Heymann *et al.*, "Beamnet: End-to-end training of a beamformer-supported multi-channel ASR system," in *ICASSP*, 2017.

[21] L. Drude *et al.*, "Integrating neural network based beamforming and weighted prediction error dereverberation." in *INTERSPEECH*, 2018.

[22] K. Kinoshita *et al.*, "Neural network-based spectrum estimation for online wpe dereverberation." in *INTERSPEECH*, 2017.

[23] J. Heymann *et al.*, "Joint optimization of neural network-based WPE dereverberation and acoustic model for robust online ASR," in *ICASSP*, 2019.

[24] T. Nakatani *et al.*, "Speech dereverberation based on variance-normalized delayed linear prediction," *IEEE-ACM T AUDIO SPE*, 2010.

[25] K. Furuya *et al.*, "Robust speech dereverberation using multichannel blind deconvolution with spectral subtraction," *IEEE-ACM T AUDIO SPE*, 2007.

[26] T. Nakatani *et al.*, "Blind speech dereverberation with multi-channel linear prediction based on short time fourier transform representation," in *ICASSP*, 2008.

[27] G. Huang *et al.*, "Kronecker product multichannel linear filtering for adaptive weighted prediction error-based speech dereverberation," *IEEE-ACM T AUDIO SPE*, 2022.

[28] D. S. Williamson and D. Wang, "Time-frequency masking in the complex domain for speech dereverberation and denoising," *IEEE-ACM T AUDIO SPE*, 2017.

[29] Y. Fu *et al.*, "Uformer: A unet based dilated complex & real dual-path conformer network for simultaneous speech enhancement and dereverberation," in *ICASSP*, 2022.

[30] Z. Wang *et al.*, "Multi-microphone complex spectral mapping for speech dereverberation," in *ICASSP*, 2020.

[31] Y. Zhao *et al.*, "Monaural speech dereverberation using temporal convolutional networks with self attention," *IEEE-ACM T AUDIO SPE*, 2020.

[32] D. Michelsanti *et al.*, "An overview of deep-learning-based audio-visual speech enhancement and separation," *IEEE-ACM T AUDIO SPE*, 2021.

[33] T. Afouras *et al.*, "The conversation: Deep audio-visual speech enhancement," in *INTERSPEECH*, 2018.

[34] J. Wu *et al.*, "Time domain audio visual speech separation," in *ASRU*, 2019.

[35] A. Ephrat *et al.*, "Looking to listen at the cocktail party: a speaker-independent audio-visual model for speech separation," *ACM T GRAPHIC*, 2018.

[36] R. Gu *et al.*, "Multi-modal multi-channel target speech separation," *IEEE J-STSP*, 2020.

[37] H. Sato *et al.*, "Multimodal attention fusion for target speaker extraction," in *SLT*, 2021.

[38] T. Ochiai *et al.*, "Multimodal speakerbeam: Single channel target speech extraction with audio-visual speaker clues." in *INTERSPEECH*, 2019.

[39] Y. Xu *et al.*, "Neural spatio-temporal beamformer for target speech separation," in *INTERSPEECH*, 2020.

[40] Y. Wu *et al.*, "Time-domain audio-visual speech separation on low quality videos," in *ICASSP*, 2022.

[41] C. Li *et al.*, "Deep audio-visual speech separation with attention mechanism," in *ICASSP*, 2020.

[42] R. Lu, Z. Duan, and C. Zhang, "Audio–visual deep clustering for speech separation," *IEEE-ACM T AUDIO SPE*, 2019.

[43] S.-W. Chung *et al.*, "Facefilter: Audio-visual speech separation using still images," in *INTERSPEECH*, 2020.

[44] W. Wang *et al.*, "A robust audio-visual speech enhancement model," in *ICASSP*, 2020.

[45] M. L. Iuzzolino *et al.*, "Av (se) 2: Audio-visual squeeze-excite speech enhancement," in *ICASSP*, 2020.

[46] G. Morrone *et al.*, "Face landmark-based speaker-independent audio-visual speech enhancement in multi-talker environments," in *ICASSP*, 2019.

[47] A. Gabbay *et al.*, "Seeing through noise: Visually driven speaker separation and enhancement," in *ICASSP*, 2018.

[48] K. Tan *et al.*, "Audio-visual speech separation and dereverberation with a two-stage multimodal network," *IEEE J-STSP*, 2020.

[49] D. Michelsanti *et al.*, "On training targets and objective functions for deep-learning-based audio-visual speech enhancement," in *ICASSP*, 2019.

[50] T. Afouras *et al.*, "Deep audio-visual speech recognition," *IEEE T PATTERN ANAL*, 2018.

[51] K. Noda *et al.*, "Audio-visual speech recognition using deep learning," *APPL INTELL*, 2015.

[52] Y. Mroueh *et al.*, "Deep multimodal learning for audio-visual speech recognition," in *ICASSP*, 2015.

[53] S. Zhang *et al.*, "Robust audio-visual speech recognition using bimodal DFSMN with multi-condition training and dropout regularization," in *ICASSP*, 2019.

[54] Y. Wu *et al.*, "Audio-visual multi-talker speech recognition in a cocktail party." in *INTERSPEECH*, 2021.

[55] S. Petridis *et al.*, "Audio-visual speech recognition with a hybrid ctc/attention architecture," in *SLT*, 2018.

[56] F. Tao *et al.*, "Gating neural network for large vocabulary audiovisual speech recognition," *IEEE-ACM T AUDIO SPE*, 2018.

[57] R. Su *et al.*, "Cross-domain deep visual feature generation for mandarin audio–visual speech recognition," *IEEE-ACM T AUDIO SPE*, 2019.

[58] A. H. Abdelaziz, "Comparing fusion models for dnn-based audiovisual continuous speech recognition," *IEEE-ACM T AUDIO SPE*, 2017.

[59] B. Shi *et al.*, "Robust self-supervised audio-visual speech recognition," in *INTERSPEECH*, 2022.

[60] P. Ma *et al.*, "End-to-end audio-visual speech recognition with conformers," in *ICASSP*, 2021.

[61] S. Petridis *et al.*, "End-to-end audiovisual speech recognition," in *ICASSP*, 2018.

[62] O. Braga *et al.*, "End-to-end multi-person audio/visual automatic speech recognition," in *ICASSP*, 2020.

[63] W. Yu *et al.*, "Fusing information streams in end-to-end audio-visual speech recognition," in *ICASSP*, 2021.

[64] W. Li *et al.*, "Improving audio-visual speech recognition performance with cross-modal student-teacher training," in *ICASSP*, 2019.

[65] P. Zhou *et al.*, "Modality attention for end-to-end audio-visual speech recognition," in *ICASSP*, 2019.

[66] R. Rose *et al.*, "End-to-end multi-talker audio-visual asr using an active speaker attention module," in *INTERSPEECH*, 2022.

[67] J. Yu *et al.*, "Audio-visual multi-channel integration and recognition of overlapped speech," *IEEE-ACM T AUDIO SPE*, 2021.

[68] G. Li *et al.*, "Audio-visual multi-channel speech separation, dereverberation and recognition," in *ICASSP*, 2022.

[69] T. Nakatani and K. Kinoshita, "A unified convolutional beamformer for simultaneous denoising and dereverberation," *IEEE Signal Process. Lett.*, 2019.

[70] T. Nakatani *et al.*, "Simultaneous denoising and dereverberation for low-latency applications using frame-by-frame online unified convolutional beamformer." in *INTERSPEECH*, 2019.

[71] T. Nakatani *et al.*, "Jointly optimal denoising, dereverberation, and source separation," *IEEE-ACM T AUDIO SPE*, vol. 28, 2020.

[72] W. Zhang *et al.*, "End-to-end dereverberation, beamforming, and speech recognition in a cocktail party," *IEEE-ACM T AUDIO SPE*, 2022.

[73] Z. Ni *et al.*, "Wpd++: An improved neural beamformer for simultaneous speech separation and dereverberation," in *SLT*, 2021.

[74] M. Delcroix *et al.*, "Strategies for distant speech recognition in reverberant environments," *EURASIP J Adv Signal Process*, 2015.

[75] T. Nakatani *et al.*, "Dnn-supported mask-based convolutional beamforming for simultaneous denoising, dereverberation, and source separation," in *ICASSP*, 2020.

[76] C. Boeddeker *et al.*, "Jointly optimal dereverberation and beamforming," in *ICASSP*, 2020.

[77] T. Nakatani *et al.*, "Blind and neural network-guided convolutional beamformer for joint denoising, dereverberation, and source separation," in *ICASSP*, 2021.

[78] T. von Neumann *et al.*, "Multi-talker asr for an unknown number of sources: Joint training of source counting, separation and asr," in *INTERSPEECH*, 2020.

[79] W. Wang *et al.*, "The sjtu system for multimodal information based speech processing challenge 2021," in *ICASSP*, 2022.

[80] A. S. Subramanian *et al.*, "Far-field location guided target speech extraction using end-to-end speech recognition objectives," in *ICASSP*, 2020.

[81] Y. Shao *et al.*, "Multi-channel multi-speaker asr using 3d spatial feature," in *ICASSP*, 2022.

[82] R. Kumar *et al.*, "End-to-end speech recognition with joint dereverberation of sub-band autoregressive envelopes," in *ICASSP*, 2022.

[83] S. Watanabe *et al.*, "Hybrid ctc/attention architecture for end-to-end speech recognition," *IEEE J-STSP*, 2017.

[84] J. Son Chung *et al.*, "Lip reading sentences in the wild," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017.

[85] I.-T. Recommendation, "Perceptual evaluation of speech quality (PESQ): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs," *Rec. ITU-T P. 862*, 2001.

[86] C. H. Taal *et al.*, "An algorithm for intelligibility prediction of time–frequency weighted noisy speech," *IEEE-ACM T AUDIO SPE*, 2011.

[87] T. H. Falk *et al.*, "A non-intrusive quality and intelligibility measure of reverberant and dereverberated speech," *IEEE-ACM T AUDIO SPE*, 2010.

[88] L. Gillick and S. J. Cox, "Some statistical issues in the comparison of speech recognition algorithms," in *ICASSP*, 1989.

[89] Z. Chen *et al.*, "Multi-channel overlapped speech recognition with location guided speech extraction network," in *SLT*, 2018.

[90] Y. Luo and N. Mesgarani, "Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation," *IEEE-ACM T AUDIO SPE*, 2019.

[91] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017.

[92] K. He *et al.*, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *IEEE Int. Conf. Comput. Vis.*, 2015.

[93] T. Afouras *et al.*, "Deep lip reading: a comparison of models and an online application," in *INTERSPEECH*, 2018.

[94] K. He *et al.*, "Deep residual learning for image recognition," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016.

[95] A. Gulati *et al.*, "Conformer: Convolution-augmented transformer for speech recognition," in *INTERSPEECH*, 2020.

[96] P. Guo *et al.*, "Recent developments on espnet toolkit boosted by conformer," in *ICASSP*, 2021.

[97] T. Ochiai *et al.*, "Beam-tasnet: Time-domain audio separation network meets frequency-domain beamformer," in *ICASSP*, 2020.

[98] H. Erdogan *et al.*, "Improved MVDR beamforming using single-channel mask prediction networks." in *INTERSPEECH*, 2016.

[99] T. Ko *et al.*, "A study on data augmentation of reverberant speech for robust speech recognition," in *ICASSP*, 2017.

[100] E. A. Habets, "Room impulse response generator," *Technische Universiteit Eindhoven, Tech. Rep*, 2006.

[101] W. Smith *et al.*, "A note on the inversion of complex matrices," *IEEE Trans. Automat. Contr.*, 1974.

[102] T. Afouras *et al.*, "Lrs3-ted: a large-scale dataset for visual speech recognition," *arXiv preprint arXiv:1809.00496*, 2018.