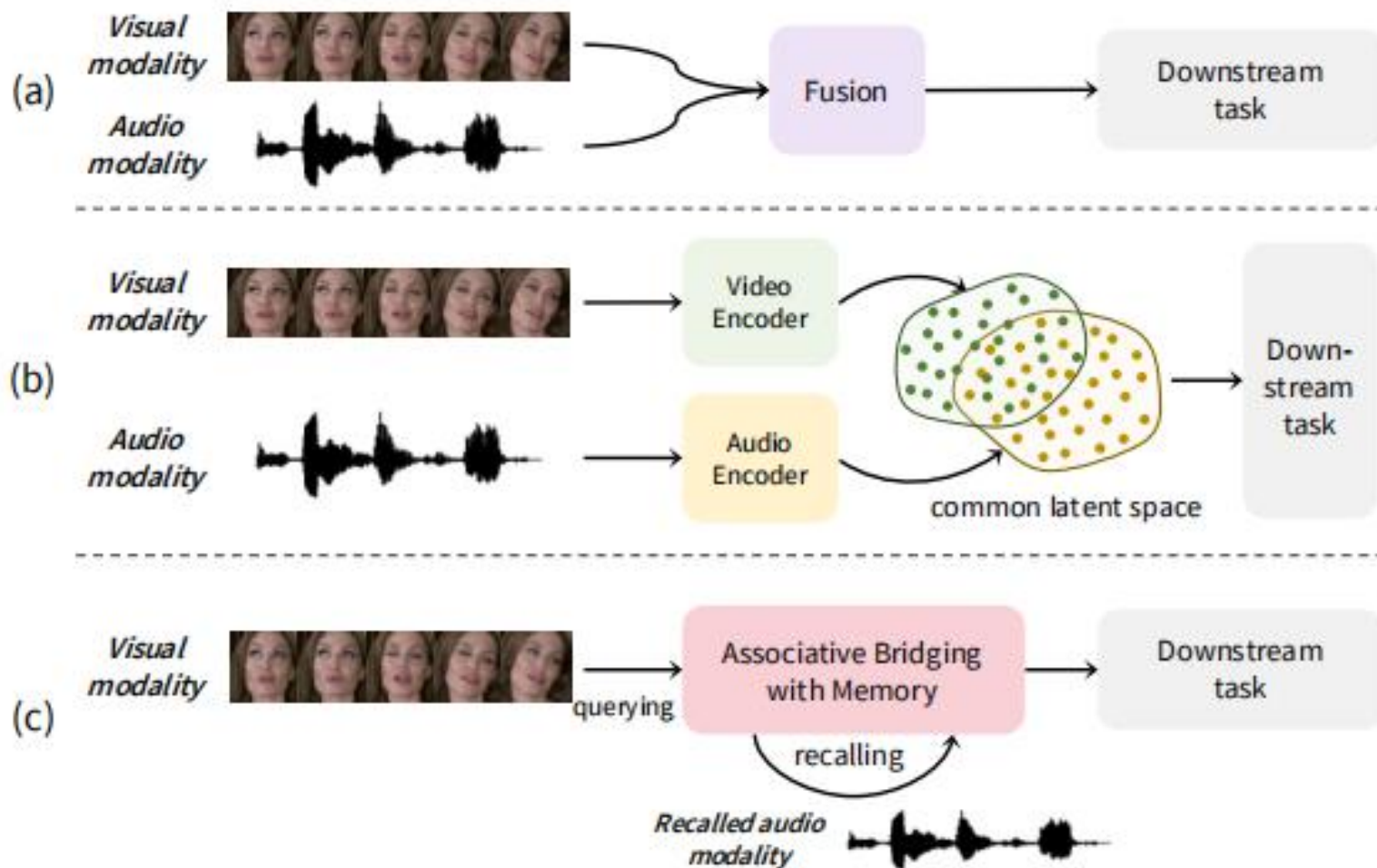


# **Multi-modality Associative Bridging through Memory: Speech Sound Recollected from Face Video**

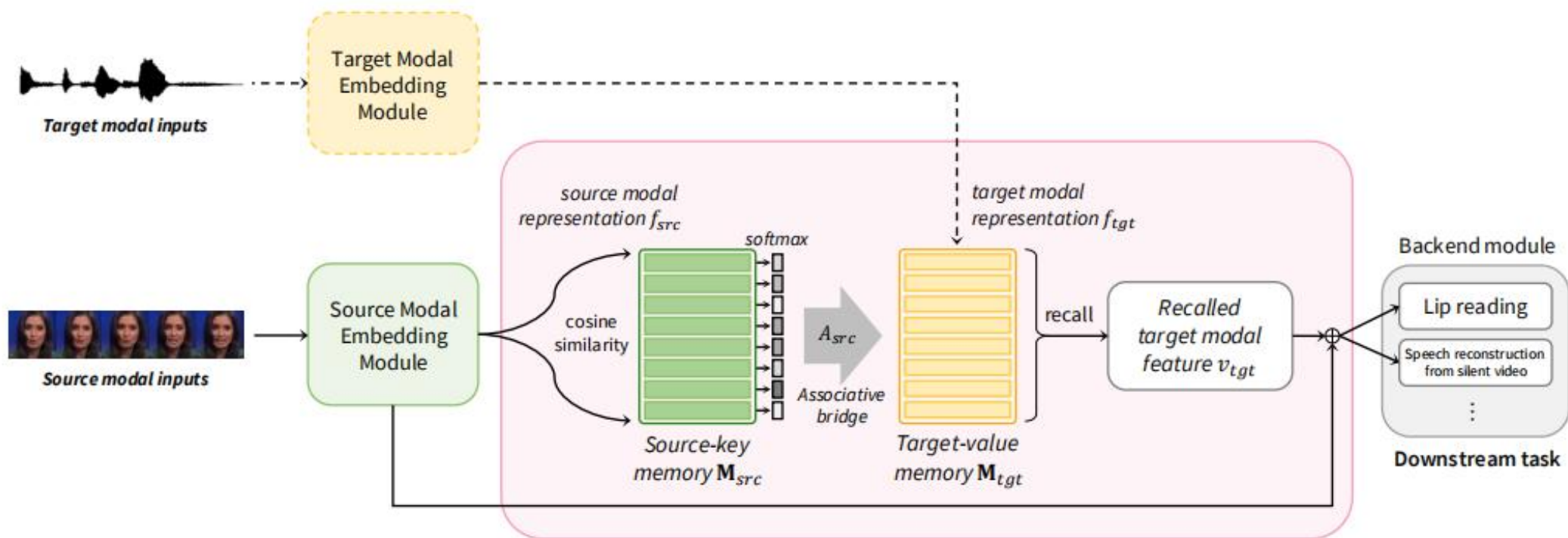
江昊宇  
2022/1/19

[https://openaccess.thecvf.com/content/ICCV2021/papers/  
Kim\\_Multi-Modality\\_Associative\\_Bridging\\_Through\\_Memory\\_Speech\\_Sound\\_Recollected\\_From\\_Face\\_ICCV\\_2021\\_paper.pdf](https://openaccess.thecvf.com/content/ICCV2021/papers/Kim_Multi-Modality_Associative_Bridging_Through_Memory_Speech_Sound_Recollected_From_Face_ICCV_2021_paper.pdf)

# Audio-Visual multi-modal learning



# Multi-modality Associative Bridging



---- : only for training

Multi-modality Associative Bridging with Memory

# Multi-modality Associative Bridging

$$s_{src}^{i,j} = \frac{\mathbf{M}_{src}^i \cdot f_{src}^j}{\|\mathbf{M}_{src}^i\|_2 \cdot \|f_{src}^j\|_2} \quad \hat{f}_{tgt}^j = A_{tgt}^j \cdot \mathbf{M}_{tgt}$$

$$\alpha_{src}^{i,j} = \frac{\exp(r \cdot s_{src}^{i,j})}{\sum_{k=1}^N \exp(r \cdot s_{src}^{k,j})} \quad v_{tgt}^j = A_{src}^j \cdot \mathbf{M}_{tgt}$$

$$A_{src}^j = \left\{ \alpha_{src}^{1,j}, \alpha_{src}^{2,j}, \dots, \alpha_{src}^{2,j}, \dots, \alpha_{tgt}^{N,j} \right\}$$

$$\mathcal{L}_{save} = \mathbb{E}_j \left[ \left\| f_{tgt}^j - \hat{f}_{tgt}^j \right\|_2^2 \right]$$

$$\mathcal{L}_{bridge} = \mathbb{E}_j \left[ D_{KL} \left( A_{tgt}^j \| A_{src}^j \right) \right]$$

$$\mathcal{L}_{task} = g(h(f_{src} \oplus v_{tgt}); y) + g(h(f_{src} \oplus f_{tgt}); y)$$

$$\mathcal{L}_{total} = \mathcal{L}_{save} + \mathcal{L}_{bridge} + \mathcal{L}_{task}$$

---

**Algorithm 1** Training algorithm of the proposed framework

---

1: **Inputs:** The training pairs of source and target modal inputs  $(X_{src}, X_{tgt})$  and label  $y$ , where  $X_{src} = \{x_{src}^l\}_{l=1}^L$ ,  $X_{tgt} = \{x_{tgt}^s\}_{s=1}^S$ . The learning rate  $\eta$ .

2: **Output:** The optimized parameters of the network  $\Phi$

---

3: Randomly initialize parameters of the network  $\Phi$

4: **for** each iteration **do**

5:  $f_{src} = \{f_{src}^j\}_{j=1}^T = \text{Source\_embed}(X_{src})$

6:  $f_{tgt} = \{f_{tgt}^j\}_{j=1}^T = \text{Target\_embed}(X_{tgt})$

7: **for**  $j = 1, 2, \dots, T$  **do**

8:  $A_{src}^j = \text{Softmax}(r \cdot \text{CosineSim}(\mathbf{M}_{src}, f_{src}^j))$

9:  $A_{tgt}^j = \text{Softmax}(r \cdot \text{CosineSim}(\mathbf{M}_{tgt}, f_{tgt}^j))$

10:  $\hat{f}_{tgt}^j = A_{tgt}^j \cdot \mathbf{M}_{tgt}$

11:  $v_{tgt}^j = A_{src}^j \cdot \mathbf{M}_{tgt}$

12: **end for**

13:  $\mathcal{L}_{save} = \sum_{j=1}^T \|f_{tgt}^j - \hat{f}_{tgt}^j\|_2^2$

14:  $\mathcal{L}_{bridge} = \sum_{j=1}^T D_{KL}(A_{tgt}^j \| A_{src}^j)$

15:  $\mathcal{L}_{task} = g(h(f_{src} \oplus v_{tgt}); y) + g(h(f_{src} \oplus f_{tgt}); y)$

16:  $\mathcal{L}_{tot} = \mathcal{L}_{save}/T + \mathcal{L}_{bridge}/T + \mathcal{L}_{task}$

17: Update  $\Phi \leftarrow \Phi - \eta \nabla_{\Phi} \mathcal{L}_{tot}$

18: **end for**

---

# Experiments Results

- Lip reading

Method	LRW	LRW-1000
Yang <i>et al.</i> [60]	83.0	38.19
Multi-Grained [51]	83.3	36.91
PCPG [30]	83.5	38.70
Deformation Flow [56]	84.1	41.93
MI Maximization [62]	84.4	38.79
Face Cutout [61]	85.0	45.24
MS-TCN [31]	85.3	41.40
<b>Proposed Method</b>	<b>85.4</b>	<b>50.82</b>

Table 1. Lip reading word accuracy comparison with visual modal inputs on LRW and LRW-1000 dataset.

N=0, 44, 88, 132 for English and N=0, 56, 112, 168 for Mandarin.

For LRW, the best word accuracy of 85.41% is achieved when N=88. The proposed framework improves the baseline with a margin of 1.27%. For LRW-1000, the best word accuracy is 50.82% when N=112 by improving the baseline performance with 5.89%. The proposed framework improves the performance regardless of the number of memory slots from the baseline in both languages.

# Experiments Results

- Speech reconstruction from silent video

Method	STOI	ESTOI	PESQ
Vid2Speech [13]	0.491	0.335	1.734
Lip2AudSpec [3]	0.513	0.352	1.673
Vougioukas <i>et al.</i> [50]	0.564	0.361	1.684
Ephrat <i>et al.</i> [12]	0.659	0.376	1.825
Lip2Wav [39]	0.731	0.535	1.772
Yadav <i>et al.</i> [58]	0.724	0.540	1.932
<b>Proposed Method</b>	<b>0.738</b>	<b>0.579</b>	<b>1.984</b>

Table 2. Performance of speech reconstruction comparison with visual modal inputs in a speaker-dependent setting on GRID.

Method	Naturalness	Intelligibility
Vid2Speech [13]	1.31 $\pm$ 0.24	1.42 $\pm$ 0.23
Lip2Wav [39]	2.83 $\pm$ 0.21	2.94 $\pm$ 0.19
<b>Proposed Method</b>	<b>2.93 <math>\pm</math>0.21</b>	<b>3.56 <math>\pm</math>0.19</b>
<b>Proposed Method</b> (+WaveNet vocoder [59])	4.37 $\pm$ 0.16	4.27 $\pm$ 0.14
Ground Truth	4.62 $\pm$ 0.13	4.57 $\pm$ 0.14

Table 3. Mean opinion scores for human evaluation on GRID.

Moreover, with WaveNet vocoder instead of Griffin-Lim, we can improve the scores as close to that of the ground truth.

# Experiments Results

- Speech reconstruction from silent video

Method	STOI	ESTOI	PESQ
Vougioukas <i>et al.</i> [50]	0.445	-	1.240
Lip2Wav [39]	0.565	0.279	1.279
<b>Proposed Method</b>	<b>0.600</b>	<b>0.315</b>	<b>1.332</b>

Table 4. Performance of speech reconstruction comparison with visual modal inputs on the speaker-independent setting on GRID.

We conduct an ablation study on different memory slot size, which is shown in supplementary material. It shows the best scores of 0.738 STOI, 0.579 ESTOI, and 1.984 PESQ when N=150. Moreover, the performance of the proposed framework improves regardless of the number of memory slots, which verifies its effectiveness.

# Experiments Results

- Learned representation inside memory

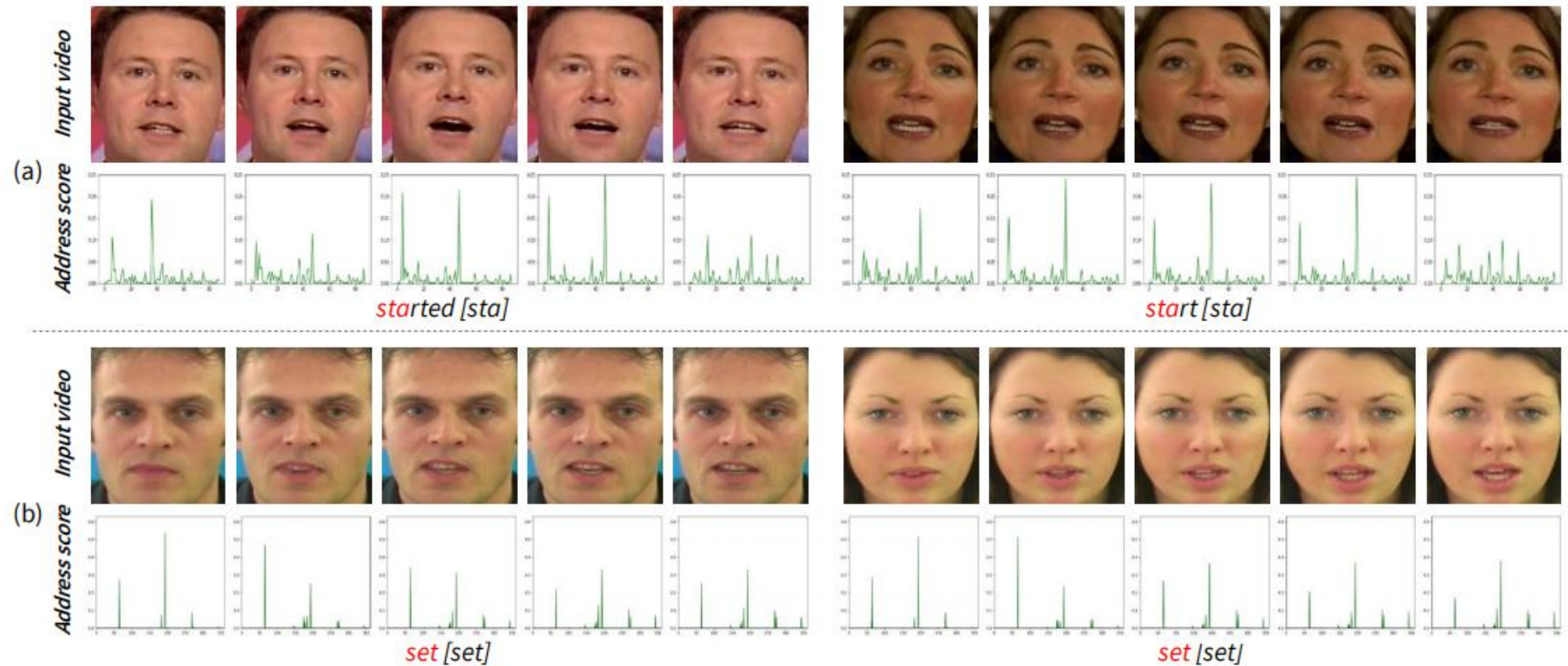


Figure 3. Face video clips (source modality) and corresponding addressing vectors for recalling audio modality (target modality) from learned representations inside memory: (a) results from lip reading and (b) results from speech reconstruction from silent video.



# Experiments Results

- Learned representation inside memory

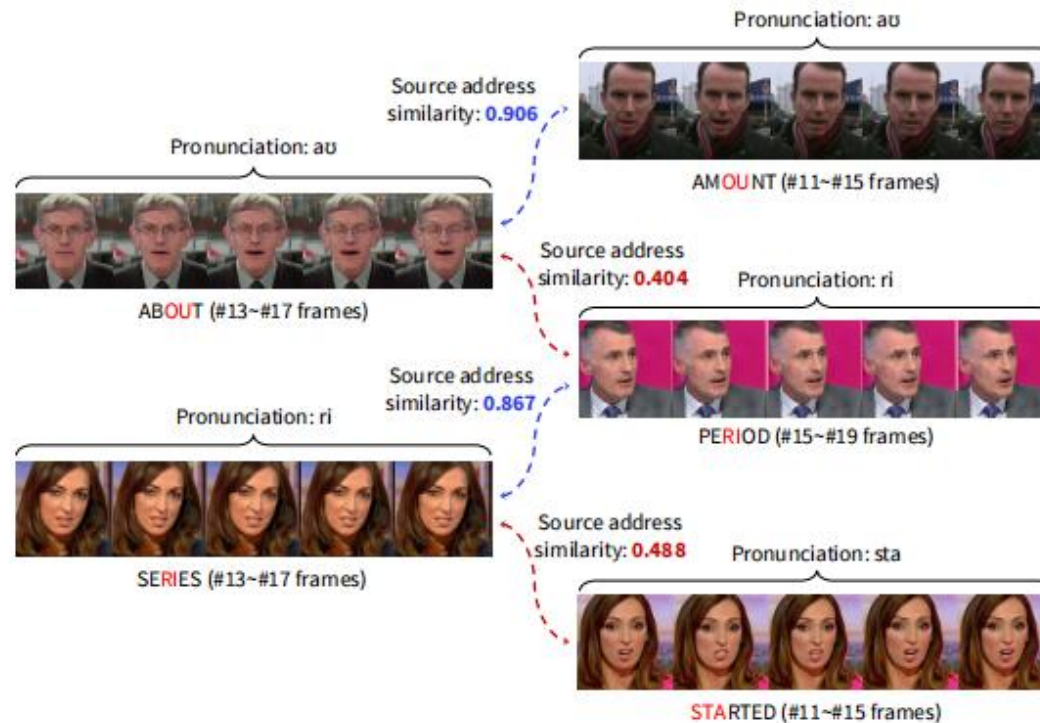


Figure 4. Examples of similarity between memory addressing vectors of different video clips in LRW. Note source addressing vector is for bridging video and audio modal features in memory.

# Experiments Results

- Comparison with methods finding a common latent space of multi-modality

Method	Baseline	Cross-modal Adaptation [7]	Knowledge Distillation [18]	<b>Proposed Method</b>
ACC(%)	84.14	84.20	84.50	<b>85.41</b>

Table 5. Lip reading word accuracy comparison with learning methods of finding a common representation of multi-modality.

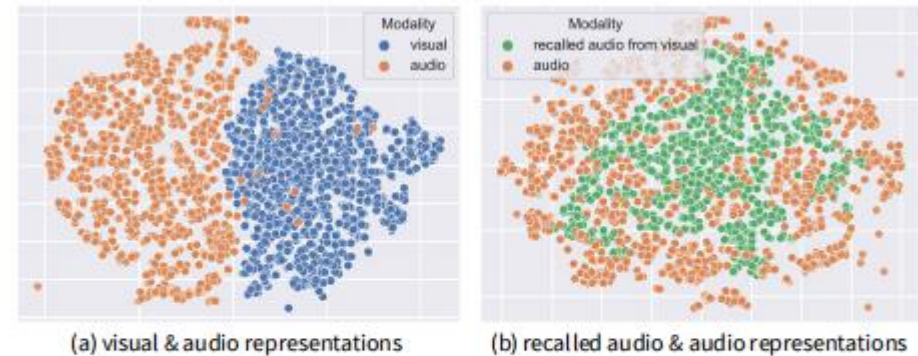


Figure 5. t-SNE [49] visualization of learned representation of (a) visual and audio modality, and (b) the recalled audio from visual modality and the actual audio modality.

# Conclusion

- With this audio-visual multi-modal bridging framework, that can utilize both audio and visual information, even with uni-modal inputs.
- The proposed framework achieves the most advanced performance in both lip-reading and speech reconstruction from silent video.
- the proposed framework can bypass the difficulty of finding a common representation of different modalities while bridging them.

Thanks