

CROSS-MODAL MUTUAL LEARNING FOR CUED SPEECH RECOGNITION

Lei Liu¹, Li Liu^{2,*}

¹Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen

²The Hong Kong University of Science and Technology (Guangzhou)

ABSTRACT

Automatic Cued Speech Recognition (ACSR) provides an intelligent human-machine interface for visual communications, where the Cued Speech (CS) system utilizes lip movements and hand gestures to code spoken language for hearing-impaired people. Previous ACSR approaches often utilize direct feature concatenation as the main fusion paradigm. However, the asynchronous modalities (*i.e.*, lip, hand shape and hand position) in CS may cause interference for feature concatenation. To address this challenge, we propose a transformer based cross-modal mutual learning framework to prompt multi-modal interaction. Compared with the vanilla self-attention, our model forces modality-specific information of different modalities to pass through a modality-invariant codebook, concatenating linguistic representations with tokens of each modality. Then the shared linguistic knowledge is used to re-synchronize multi-modal sequences. Moreover, we establish a novel large-scale multi-speaker CS dataset for Mandarin Chinese. To our knowledge, this is the first work on ACSR for Mandarin Chinese. Extensive experiments are conducted for different languages (*i.e.*, Chinese, French, and British English). Results demonstrate that our model exhibits superior recognition performance to the state-of-the-art by a large margin.

Index Terms— Mandarin Chinese Cued Speech, Multi-modal Transformer, Linguistic Representation

1. INTRODUCTION

Cued Speech (CS) [1, 2] is an efficient communication system for hearing impaired people, which leverages lip motions and hand gestures into visual cues. As shown in Figure 1 in Mandarin Chinese CS [3, 4], hand shapes and positions are the supplementary to alleviate the visual ambiguity of similar labial shapes caused by lip reading (*e.g.*, [p] and [b]). Taking lip and hand as two distinct modalities, Automatic Cued Speech Recognition (ACSR) aims to recognize the multi-modal inputs into linguistic text. The main challenge for ACSR is the natural asynchrony between lip-hand movements, *i.e.*, the hand (*i.e.*, hand shape and hand position)

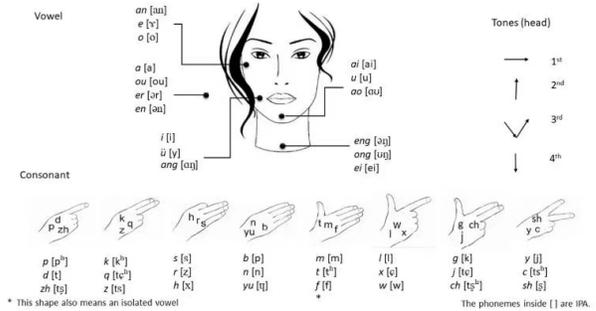


Fig. 1. The Mandarin Chinese CS system (from [3]).

generally move faster than lips to prepare the next phoneme in a CS system, called hand preceding phenomenon [5, 6].

Existing researches mainly assumed lip-hand movements are synchronous by default, thus tended to extract discriminative representations and directly concatenate the high-level features of multi-modal inputs. For example, [7, 8] utilized artificial marks to obtain regions of interests (ROIs) and directly concatenated features of lip and hand. MSHMM [9] merged different features by giving weights for different CS modalities. These methods simply exploited feature concatenation as the fusion strategy while ignoring the asynchronous issue, as well as recent knowledge distillation based method [10]. To handle asynchronous modalities in the ACSR task, a re-synchronization procedure [6] was proposed to align visual features using the prior of hand preceding time, which is derived from the statistical information of CS speakers. However, such prior is dependent on speakers and datasets, which is difficult to generalize to different languages. Therefore, there still lacks of effective and flexible fusion strategies for asynchronous CS modalities. Besides, by taking into account the context relationships of phonemes in long-time CS videos, it would be desirable to capture global dependency [11] over dynamic longer context for solving the above cross-modal fusion task, which could further enhance the interaction of multi-modal sequences in CS. Moreover, previous ACSR task generally utilized phoneme-level annotations on two small scale datasets in French and British English.

In this study, we firstly establish a novel large-scale multi-speaker CS video dataset for Mandarin Chinese. Then we

* Corresponding author: avrilliu@hkust-gz.edu.cn.

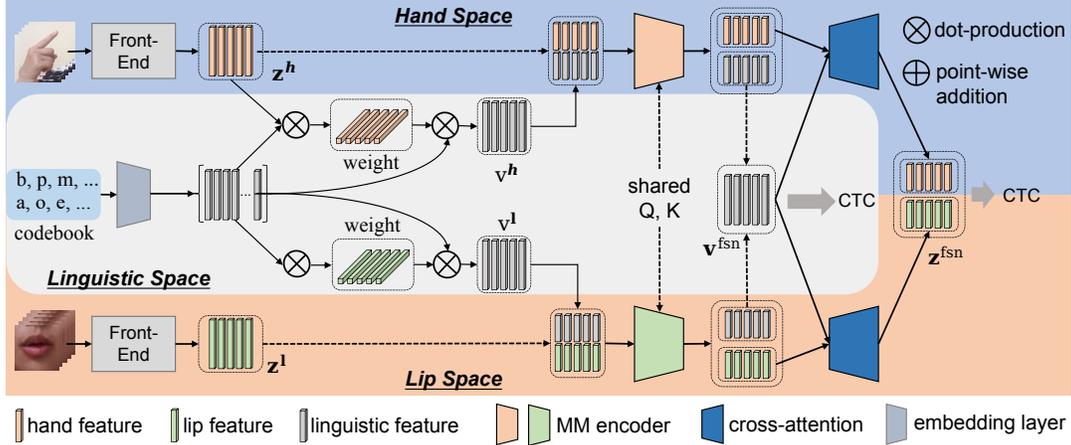


Fig. 2. The proposed framework for ACSR. A shared front-end is deployed to extract frame-wise features for each modality and a codebook is to generate modality-invariant linguistic features shared by different modalities. Then a multi-modal (MM) encoder is to capture long-time temporal dependencies, where encoders of different modalities share query and key weights for interactions. Referring to linguistic information, a cross-modal alignment is conducted for visual representations.

propose a novel transformer based cross-modal mutual learning framework (see Figure 2) to prompt cross-modal information interaction, which achieves feature fusion based on the aligned modalities with the long-time dependencies. In particular, a pre-trained front-end is used to extract frame-wise representations. Then a multi-modal encoder with sharing query and key weights is used to enhance cross-modal interactions. To capture modality-invariant information, our method forces visual representations of different modalities to pass through a linguistic codebook, requiring the model to collate a compact linguistic representation for tokens of each modality. Then linguistic knowledge is exploited to align different modalities based on the cross-attention mechanism.

In summary, the key contributions are: **(1)** We propose a novel transformer to prompt cross-modal mutual learning for ACSR task based on re-asynchronous modalities with long-time dependencies. Notably, this is the first work that only requires sentence-level annotations to handle asynchronous issue, unlike previous methods relying phoneme-level annotations. **(2)** We propose to collate a modality-invariant shared representation for multi-modal tokens to obtain linguistic information, which is utilized to guide the alignments for multi-modal data stream. **(3)** To the best of knowledge, this is the first work on ACSR for Mandarin Chinese. Extensive experiments on three CS benchmarks (*i.e.*, Chinese, French, and British English) demonstrate that our model can significantly outperform the state-of-the-art (SOTA).

2. METHODOLOGY

Problem Formulation. Given a CS dataset of N quadruples $\mathcal{D} = \{(x_i^l, x_i^g, x_i^p, y_i)\}_{i=1}^N$, corresponding to the lip, hand shape, hand position, and sentence-level label sequences.

Our target is to learn linguistic features represented by multi-modal data stream, where lip and hand sequences are complementary to each other as different modalities. Besides, previous works [6, 10] often exploited hand region of interests (ROIs) and position coordinates to extract hand shape and position features, respectively. In this work, we experimentally demonstrate that hand ROIs contain both shape and position information, thus only take hand ROIs as the hand input. The proposed framework is shown in Figure 2.

2.1. Cross-Modal Mutual Learning

Front-end. The front-end adopts a modified ResNet-18 [12] where the first layer is replaced by a 3D convolutional layer with a kernel size $5 \times 3 \times 3$. The features at the penultimate layer are squeezed along the spatial dimension by global average pooling, resulting in modality-specific features, *i.e.*, $z^l, z^g, z^p \in \mathcal{R}^d$ for lip, hand shape, and hand position, where d is the feature dimension. Element-wise addition operation \oplus is conducted to fuse features of hands via $z^h = z^g \oplus z^p$. To be simplified, we denote $m \in \{l, h\}$ as lip (l) and hand (h) modalities in the following section, respectively.

Modality-Invariant Linguistic Codebook. In this work, we aim to extract modality-invariant linguistic representation from the visual features z^l and z^h . Inspired by [13], we propose to learn a CS codebook to extract linguistic information shared by lip and hand sequences. To this end, an embedding layer is exploited to obtain linguistic codebook basis of cued speech codes. Given the codebook base set $\mathbf{D} = \{b_i\}_{i=1}^n$, where $b_i \in \mathcal{R}^d$ and n is the total number of basis. To derive linguistic information, we compute the dot products of the frame-wise visual features (z^l or z^h) with all bases of \mathbf{D} , and apply a softmax function to obtain the normalized weights.

This produces the modality-invariant linguistic representations by $\mathbf{v}^m = \text{softmax}(\mathbf{D}\mathbf{z}^m)\mathbf{D}$, where $m \in \{l, h\}$ and linguistic representation is the weighted sum of the bases and the dot product denotes pairwise similarity between \mathbf{z}^l (or \mathbf{z}^h) and each basis of \mathbf{D} .

Multi-Modal Encoder. Given the visual and linguistic representations for lip and hand stream, we aim to force the model to emphasize the temporal information in both modality-invariant and modality-specific flows. Hence, we restrict the cross-modal interactions along with the full sequence to capture long-time dependencies. Formally, we concatenate the visual and linguistic features into a single sequence as:

$$\mathbf{u}^m = [\mathbf{z}^m || \hat{\mathbf{v}}^m], \quad \text{where} \quad \hat{\mathbf{v}}^m = g(\mathbf{v}^m, \mathbf{E}_{\text{ling}}^{\text{sub}}), \quad (1)$$

where $m \in \{l, h\}$ and $[\cdot || \cdot]$ denotes the frame-wise concatenation of the tokens between \mathbf{z}^l (or \mathbf{z}^h) and $\hat{\mathbf{v}}^l$ (or $\hat{\mathbf{v}}^h$). We use the projection $\mathbf{E}_{\text{ling}}^{\text{sub}}$ to reduce the dimension of linguistic representations, which can decrease the computation complexity of pairwise attention. The multi-modal encoder stacks a series of vanilla transformer layers, where each modality has its own dedicated value parameters (*i.e.*, θ_v^l and θ_v^h), and shared query θ_q and key θ_k parameters. Each lip (hand) token can attend to all other lip (hand) and linguistic tokens via self-attention:

$$\mathbf{u}_{L+1}^m = \text{Transformer}(\mathbf{u}_L^m; \theta_q, \theta_k, \theta_v^m), \quad (2)$$

where $m \in \{l, h\}$ and L denotes L -th transformer layer with vanilla self-attention blocks. This layer allows the information exchanging between visual and linguistic tokens.

We can generalise this model by allowing to linguistic interactions between lip and hand modalities, which is achieved by sharing modality information using compact linguistic vectors. The tokens at layer L are calculated as:

$$\mathbf{z}_{L+1}^m || \hat{\mathbf{v}}_{L+1}^m = \text{Transformer}(\mathbf{z}_L^m, \hat{\mathbf{v}}_L^{\text{fsn}}; \theta_q, \theta_k, \theta_v^m), \quad (3)$$

$$\hat{\mathbf{v}}_L^{\text{fsn}} = \frac{1}{2}(\hat{\mathbf{v}}_L^l + \hat{\mathbf{v}}_L^h),$$

where $m \in \{l, h\}$. Since multi-modal attention flows must pass through the shared tight fused linguistics, \mathbf{z}^l and \mathbf{z}^h can exchange information via \mathbf{v}^{fsn} within a transformer layer. Thus the model can enhance the cross-modal interactions by using shared linguistic information for the ACSR task.

Visual-Linguistic Alignment (VLA). To alleviate the naturally asynchronous issue in CS, we aim to integrate aligned information from multiple modalities using cross-attention mechanism. Given the encoded representations of the multi-modal encoder, *i.e.*, \mathbf{z}^l , \mathbf{z}^h , and \mathbf{v}^{fsn} , the fused linguistic information $\hat{\mathbf{v}}^{\text{fsn}}$ can be the reference to align lip and hand sequences, since linguistic knowledge is shared by both modalities. We first recover the tight linguistic representation $\hat{\mathbf{v}}^{\text{fsn}}$ into the original dimension by $\mathbf{v}^{\text{fsn}} = g(\hat{\mathbf{v}}^{\text{fsn}}, \mathbf{E}_{\text{ling}}^{\text{up}})$, where $\mathbf{E}_{\text{ling}}^{\text{up}}$ is a linear projection. Then a cross-attention layer is defined to achieve alignments as:

$$\mathbf{z}_{L+1}^m = \text{Cross-Transformer}(\mathbf{z}_L^m, \mathbf{v}^{\text{fsn}}; \theta_{q,k,v}). \quad (4)$$

Here, the vanilla cross-attention operation [11] is employed with the shared query, key, and value weights for lip and hand, allowing the multi-modal streams to align with modality-invariant linguistic information. The final fused representation \mathbf{z}^{fsn} is obtained by the frame-wise concatenation of \mathbf{z}^l and \mathbf{z}^h as $\mathbf{z}^{\text{fsn}} = [\mathbf{z}^l || \mathbf{z}^h]$.

Objective Function. Connectionist Temporal Classification (CTC) loss [14] shows superior performance for speech recognition tasks. Given the input \mathbf{x} and target sequence \mathbf{y} , CTC computes the negative log likelihood of the posterior by summing the probabilities over valid alignment set $\mathcal{A}_{\mathbf{x},\mathbf{y}}$:

$$-\ln p(\mathbf{y} | \mathbf{x}) = -\ln \sum_{A \in \mathcal{A}_{\mathbf{x},\mathbf{y}}} \prod_{t=1}^T p_t(a_t | \mathbf{x}), \quad (5)$$

where T is the target length. In this work, we utilize a hybrid CTC architecture to force prediction consistency between visual and linguistic representations:

$$\mathcal{L}_{\text{hyb}} = -\ln p(\mathbf{y} | \mathbf{z}^{\text{fsn}}) - \ln p(\mathbf{y} | \mathbf{v}^{\text{fsn}}), \quad (6)$$

where visual and linguistic representation corresponds to the same target sequence.

3. EXPERIMENTS

Datasets. Three datasets are used to evaluate the performance of the proposed method, including French CS [2], British English CS [15], and the newly collected Mandarin Chinese CS. The details of public CS datasets can refer to Table 1.

Implementations. The transformer is randomly initialized while the front-end is pretrained on ImageNet. The multi-modal encoder uses 3 multi-modal blocks. The mini-batch size is set as 1 [16]. Following [11], the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 0.05$ is used. The learning rate increases linearly with the first 4000 steps, yielding a peak learning rate and then decreases proportionally to the inverse square root of the step number. The whole network is trained for 50 epochs.

Protocols. The training and test sentences are randomly split as 4 : 1 as shown in Table 1. All methods are evaluated in character error rate (CER) and word error rate (WER) to evaluate the recognition ability on phoneme and word levels.

3.1. Experimental Results

Chinese CS Dataset. As shown in Table 2, our method achieves SOTA results compared with baseline, *i.e.*, 9.7% CER in the single-speaker setting and 24.5% CER in the multi-speaker setting. JLF [10] is the previous SOTA method using LSTM, which performs worse than self-attention based approaches. Previous methods still suffer from the asynchronous issue in Chinese CS. Our method can further improve the recognition accuracy via multi-modal alignment.

Table 1. Details of CS datasets with different languages.

| Dataset | French | British | | Chinese | |
|-----------|-----------|---------|----------|-----------|-------------|
| speaker | 1 | 1 | 5 | 1 | 4 |
| sentence | 238 | 97 | 390 | 1000 | 4000 |
| character | 12872 | 2741 | 11021 | 32902 | 131581 |
| word | - | - | - | 10562 | 42248 |
| phoneme | 35 | 44 | 44 | 40 | 40 |
| gesture | 8 | 8 | 8 | 8 | 8 |
| position | 5 | 4 | 4 | 5 | 5 |
| train | 193/10636 | 78/2240 | 312/8924 | 800/26683 | 3200/105372 |
| test | 45/2236 | 19/501 | 78/2097 | 200/6219 | 800/26209 |

¹ #Train/#Test is in the form of word/character.

Table 2. Performance comparisons on Chinese CS dataset.

| Dataset | Chinese | | | |
|----------------------|------------|-------------|-------------|-------------|
| #Speaker | single | | multiple | |
| Metrics | CER | WER | CER | WER |
| CNN + LSTM [17] | 55.4 | 92.8 | 61.4 | 96.1 |
| CNN + CTC [12] | 35.6 | 78.3 | 41.9 | 83.4 |
| JLF + COS + CTC [10] | 33.5 | 67.1 | 68.2 | 98.1 |
| Self-attention [11] | 26.1 | 61.8 | 38.8 | 78.6 |
| Ours | 9.7 | 24.1 | 24.5 | 54.5 |

Table 3. Performance comparisons (CER) on British and French CS datasets. WER is unavailable due to lacking of word-level annotations. JLF3 is the previous SOTA.

| Dataset | French | British | |
|------------------|-------------|-------------|-------------|
| #Speaker | single | single | multiple |
| CNN-HMM [9] | 38.0 | - | - |
| Fully Conv [17] | 29.2 | 36.3 | - |
| CNN + LSTM [17] | 33.4 | 43.6 | - |
| Transformer [17] | 37.5 | 39.8 | - |
| Student CE [10] | 35.6 | 47.5 | 37.7 |
| JLF1 [10] | 27.5 | 38.5 | 34.3 |
| JLF2 [10] | 27.5 | 36.9 | 31.5 |
| JLF3 [10] | 25.8 | 35.1 | 30.3 |
| Ours | 24.9 | 33.6 | 29.2 |

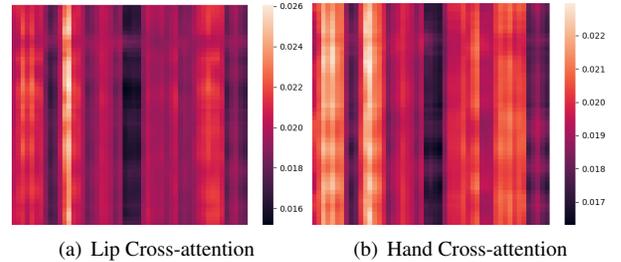
We notice that multi-speaker task is more difficult than single-speaker one due to domain adaption issue.

French&British CS Dataset. As shown in Table 3, our method can achieve best results on both French and British CS datasets, outperforming previous SOTA. Compared with LSTM and vanilla transformer, our method benefits from the aligned modalities and can capture effective long-time dependency via multi-modal interaction. The performance improvement is small due to the small data scale of these datasets. The experimental results indicate the effectiveness of the proposed method on small scale datasets.

Cross-attention Score. The cross-attention scores of the VLA module are shown in Figure 3. It is observed that lip and hand modalities exhibit similar cross-attention scores,

Table 4. Ablation Studies. CB is codebook and concat is the concentration of lip and hand features.

| Dataset | Chinese | | | |
|-----------------------------|------------|-------------|-------------|-------------|
| #Speaker | single | | multiple | |
| Metrics (100%) | CER | WER | CER | WER |
| CNN + concat | 35.6 | 78.3 | 41.9 | 83.4 |
| + self-attn | 26.1 | 61.8 | 38.8 | 78.6 |
| + self-attn + CB (Eq. 2, 3) | 23.1 | 59.0 | 36.9 | 76.7 |
| + self-attn + CB (Eq. 4) | 16.3 | 43.1 | 34.7 | 72.5 |
| + self-attn + VLA | 14.7 | 38.9 | 31.6 | 70.3 |
| + self-attn + CB + VLA | 9.7 | 24.1 | 24.5 | 54.5 |

**Fig. 3.** Cross-attention scores on Chinese CS dataset.

indicating VLA can achieve better cross-modal alignments. **Ablation Studies.** Table 4 exhibits the ablation studies for the proposed method. ‘+ self-attn + VLA’ only uses linguistic features in the VLA, while lip/hand encoder is without cross-modal interaction. Both codebook and VLA can further improve the performance, outperforming the baselines without codebook and VLA. There exists a significant performance drop when removing each component for our method, indicating the effectiveness of the proposed cross-modal strategy.

4. CONCLUSIONS

In this work, we proposed a cross-modal mutual learning framework for the ACSR task. We present a multi-modal transformer to capture long-time dependencies for both lip and hand sequences, which transforms modality-specific information into modality-invariant linguistic features via a linguistic codebook. Then modality-invariant linguistic information can guide the cross-modal alignment via the cross-attention operation. The experimental results demonstrate that the proposed approach achieves new SOTA on the ACSR. For the future work, we will be engaged in improving the CS model by decreasing the sentence-level errors.

Acknowledgments. This work is supported by the National Natural Science Foundation of China (No. 62101351), and the GuangDong Basic and Applied Basic Research Foundation (No.2020A1515110376).

5. REFERENCES

- [1] R Orin Cornett, “Cued speech,” *American Annals of the Deaf*, pp. 3–13, 1967.
- [2] Li Liu, Gang Feng, and Denis Beuitemps, “Automatic temporal segmentation of hand movements for hand positions recognition in french cued speech,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 3061–3065.
- [3] Li Liu and Gang Feng, “A pilot study on mandarin chinese cued speech,” *American Annals of the Deaf*, vol. 164, no. 4, pp. 496–518, 2019.
- [4] Li Liu, Gang Feng, Xiaoxi Ren, and Xianping Ma, “Objective hand complexity comparison between two mandarin chinese cued speech systems,” in *International Symposium on Chinese Spoken Language Processing (ISCSLP)*, 2022, pp. 215–219.
- [5] Li Liu, Gang Feng, Denis Beuitemps, and Xiao-Ping Zhang, “A novel resynchronization procedure for hand-lips fusion applied to continuous french cued speech recognition,” in *European Signal Processing Conference (EUSIPCO)*, 2019, pp. 1–5.
- [6] Li Liu, Gang Feng, Denis Beuitemps, and Xiao-Ping Zhang, “Re-synchronization using the hand preceding model for multi-modal fusion in automatic continuous cued speech recognition,” *IEEE Transactions on Multimedia*, vol. 23, pp. 292–305, 2020.
- [7] Panikos Heracleous, Denis Beuitemps, and Norihiro Hagita, “Continuous phoneme recognition in cued speech for french,” in *European Signal Processing Conference (EUSIPCO)*, 2012, pp. 2090–2093.
- [8] Jianrong Wang, Nan Gu, Mei Yu, Xuwei Li, Qiang Fang, and Li Liu, “An attention self-supervised contrastive learning based three-stage model for hand shape feature representation in cued speech,” in *Proceedings of Interspeech*, 2021, pp. 626–630.
- [9] Li Liu, Thomas Hueber, Gang Feng, and Denis Beuitemps, “Visual recognition of continuous cued speech using a tandem cnn-hmm approach,” in *Proceedings of Interspeech*, 2018, pp. 2643–2647.
- [10] Jianrong Wang, Ziyue Tang, Xuwei Li, Mei Yu, Qiang Fang, and Li Liu, “Cross-modal knowledge distillation method for automatic cued speech recognition,” in *Proceedings of Interspeech*, 2021, p. 2986–2990.
- [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, 2017, vol. 30.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [13] Chih-Chun Yang, Wan-Cyuan Fan, Cheng-Fu Yang, and Yu-Chiang Frank Wang, “Cross-modal mutual learning for audio-visual speech recognition and manipulation,” in *Proceedings of the 36th AAAI Conference on Artificial Intelligence*, 2022, vol. 22.
- [14] Pingchuan Ma, Stavros Petridis, and Maja Pantic, “Investigating the lombard effect influence on end-to-end audio-visual speech recognition,” in *Proceedings of Interspeech*, 2019, pp. 4090–4094.
- [15] Sanjana Sankar, Denis Beuitemps, and Thomas Hueber, “Multistream neural architectures for cued speech recognition using a pre-trained visual feature extractor and constrained ctc decoding,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 8477–8481.
- [16] Pingchuan Ma, Stavros Petridis, and Maja Pantic, “End-to-end audio-visual speech recognition with conformers,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7613–7617.
- [17] Katerina Papadimitriou and Gerasimos Potamianos, “A fully convolutional sequence learning approach for cued speech recognition from videos,” in *European Signal Processing Conference (EUSIPCO)*, 2021, pp. 326–330.