

# Enhanced Self-Supervised Multi-View Representations with Modality-Missing Robustness for Audio-Visual Speech Recognition

Anonymous ICME submission

**Abstract**—Audio-Visual Speech Recognition (AVSR) leverages visual information to enhance speech understanding. However, current models often assume stable, frontal viewpoints, suffering significant performance drops with non-frontal angles or when video signals are partially missing. Our approach first employs a multi-view data generation strategy using 3D head avatar reconstruction, synthesizing viewpoint-diverse data to handle varying head poses. Then, we introduce a self-supervised multi-view representation learning paradigm that integrates a Multi-View Consistency Loss and a Representation Domain Alignment Loss, ensuring that learned embeddings are robust to viewpoint shifts and domain variations. We develop a Unified Modality Adapter (UMA) module to seamlessly revert to audio-only performance levels when visual inputs are unavailable, maintaining the benefits of audio-visual modeling under modality-missing conditions. [14] Experimental results on challenging datasets demonstrating significant improvements in recognition accuracy under diverse viewing conditions including cases where visual inputs are partially or fully missing.

## II. RELATED WORK

**Lipreading and Multi-View Representation Learning.** Lipreading, which aims to decode speech solely from visual cues of the speaker's lips, has witnessed substantial progress. Early methods relied on constrained settings and often focused on isolated words or simple phrases [16]. With the introduction of large-scale lipreading such as LRSS2 [17] and LRSS3 [13], along with increasingly sophisticated deep learning models [18]–[20], lipreading has advanced towards continuous, sentence-level recognition. Furthermore, some multi-modal fusion methods [21]–[23] have been proposed to enhance audio-visual and personalized appearance. Neural rendering, combined with parametric priors, now enables photo-realistic and controllable avatars that better preserve subtle lip movements critical for AVSR under multi-view conditions [12], [13].

**Robustness to Missing Video Modality.** Real-world audio-visual speech recognition must cope with scenarios where visual inputs are degraded, partially visible, or entirely absent [18], [24]. Early strategies attempted to combine audio and video frames, encouraging models to rely more on audio cues when the visual modality failed [5], [15]. Beyond simple dropout, current methods explore architectures and training regimes that ensure the presence of visual features never diminishes performance [11]. Novel frameworks incorporate multi-view data and modality-adaptive modules to seamlessly revert to audio-only quality levels under adverse conditions [21], [36], thereby offering robust performance across a wide range of viewing angles and video availability scenarios.

**3D Facial Avatar Reconstruction.** Early work often avoided explicit 3D structures and relied on image-based warping or implicit representations [27], [28]. However, these approaches typically struggled with significant pose shifts due to limited geometric consistency. 3D Morphable Models, such as FLAME [29], [30], introduced low-dimensional shape and expression spaces, serving as foundations for mesh-based avatars that ensure geometric accuracy and temporal stability [31]. Recent advances in 3D head reconstruction have been driven

by approaches that unify neural rendering and deep learning-based geometry estimation. Grassal et al. [12] proposed a method to generate neural head avatars from monocular RGB videos, leveraging deep neural networks to produce dense 3D geometry and personalized appearance. Neural rendering, combined with parametric priors, now enables photo-realistic and controllable avatars that better preserve subtle lip movements critical for AVSR under multi-view conditions [12], [13].

Our approach adopts self-supervised pretraining strategy to learn multi-view representations from unlabeled data [24]–[26], enabling the capture of robust cross-modal correlations and visual patterns even under challenging multi-view conditions.

**3D Facial Avatar Reconstruction.** Early work often avoided explicit 3D structures and relied on image-based warping or implicit representations [27], [28]. However, these approaches typically struggled with significant pose shifts due to limited geometric consistency. 3D Morphable Models, such as FLAME [29], [30], introduced low-dimensional shape and expression spaces, serving as foundations for mesh-based avatars that ensure geometric accuracy and temporal stability [31]. Recent advances in 3D head reconstruction have been driven

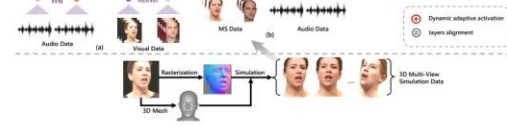


Fig. 1. Overview of the proposed framework. (a) AV-HuBERT-based Audio-Visual Encoder, computes the masked multi-modal prediction loss [15] on real data. FPN denotes a feed-forward network for audio extraction, and ResNet is a modified ResNet for video feature extraction. (b) The self-supervised Multi-View Representation Learning (MVRL) module, which leverages multi-view simulated data (MS Data) generated by the lower pipeline. Parameters in the MV Encoder are shared with the Transformer Encoder to ensure consistent representation learning. (c) The Unified Modality Adapter (UMA) module, where “AV” represents fused audio-visual features, “E” denotes UMA encoder audio embeddings and “E<sub>vis</sub>” denotes the audio-visual encoder embeddings for the same input. When UMA is activated, audio features skip the “EMA” fusion step and pass through UMA, providing robustness to missing visual inputs.

by approaches that unify neural rendering and deep learning-based geometry estimation. Grassal et al. [12] proposed a method to generate neural head avatars from monocular RGB videos, leveraging deep neural networks to produce dense 3D geometry and personalized appearance. Neural rendering, combined with parametric priors, now enables photo-realistic and controllable avatars that better preserve subtle lip movements critical for AVSR under multi-view conditions [12], [13].

**Robustness to Missing Video Modality.** Real-world audio-visual speech recognition must cope with scenarios where visual inputs are degraded, partially visible, or entirely absent [18], [24]. Early strategies attempted to combine audio and video frames, encouraging models to rely more on audio cues when the visual modality failed [5], [15]. Beyond simple dropout, current methods explore architectures and training regimes that ensure the presence of visual features never diminishes performance [11]. Novel frameworks incorporate multi-view data and modality-adaptive modules to seamlessly revert to audio-only quality levels under adverse conditions [21], [36], thereby offering robust performance across a wide range of viewing angles and video availability scenarios.

**3D Facial Avatar Reconstruction.** Early work often avoided explicit 3D structures and relied on image-based warping or implicit representations [27], [28]. However, these approaches typically struggled with significant pose shifts due to limited geometric consistency. 3D Morphable Models, such as FLAME [29], [30], introduced low-dimensional shape and expression spaces, serving as foundations for mesh-based avatars that ensure geometric accuracy and temporal stability [31]. Recent advances in 3D head reconstruction have been driven

by approaches that unify neural rendering and deep learning-based geometry estimation. Grassal et al. [12] proposed a method to generate neural head avatars from monocular RGB videos, leveraging deep neural networks to produce dense 3D geometry and personalized appearance. Neural rendering, combined with parametric priors, now enables photo-realistic and controllable avatars that better preserve subtle lip movements critical for AVSR under multi-view conditions [12], [13].

Our proposed framework is illustrated in Fig. 1, the methodology involves synthesizing multi-view data via 3D reconstruction, learning viewpoint-invariant and domain-aligned representations through self-supervised training.

## A. Multi-View Data Generation Strategy

**1) Data Preparation:** As shown in Fig. 2, we first perform head-face and face-only data generation by solving a Perspective-n-Point (PnP) problem via 2D facial landmarks, thereby extracting Euler angles [37], [38]. Subsequently, we select video segments exhibiting substantial rotational changes as the training data for multi-view data generation.

**2) Multi-View Data Simulation:** We proceed to generate additional viewing angles using a 3D head avatar reconstruction approach. Instead of relying on multiple cameras or extensive capture setups, we employ a learned neural avatar that reconstructs a detailed, animatable 3D representation of the speaker from RGB video segments [27], [28]. We adopt a 3D morphable model (e.g., FLAME [29]) as a geometric backbone. Given shape parameters  $\alpha$ , expression parameters  $\epsilon$ , and pose parameters  $\theta$ , the base model provides a coarse mesh:

$$M_{base}(\alpha, \epsilon, \theta) \in \mathbb{R}^{3 \times V \times 3}, \quad (1)$$

where  $N_V$  is the number of vertices. We use geometry refinement  $F$  to capture fine details not represented by the



Fig. 2. Multi-view data generation strategy.

template:

$$M(\alpha, \epsilon, \theta) = M_{base}(\alpha, \epsilon, \theta) + G(\theta). \quad (2)$$

The coarse geometric alignment and fine-grained texture refinement ensure lip contours align closely with the original frames. [21], [39], [40].

## B. Self-Supervised Multi-View Representation Learning

Leveraging self-supervised training to capture intrinsic audio-visual correlations [23], [25], we utilize a feature extractor and encoder structure inspired by AV-HuBERT [24] to learn viewpoint-invariant, domain-aligned audio-visual embeddings.

**1) Multi-View Consistency Loss:** Consider a batch of  $N$  paired samples, where  $\{x_{a,i}\}_{i=1}^N$  are real audio-visual sequences and  $\{x_{v,i}\}_{i=1}^N$  are their corresponding synthesized multi-view counterparts. We seek an embedding function  $f(\cdot)$  that maps these inputs into a latent space invariant to viewpoint changes. Let  $f(x_{a,i})$  and  $f(x_{v,i})$  be embeddings with dimensions  $T \times D$ , where  $T$  is the combined spatio-temporal dimension and  $D$  is the feature channel dimension.

We first consider an element-wise alignment cost to ensure that real and synthetic embeddings match at the feature level. Defining a mean squared error term:

$$L_{mse} = \frac{1}{N} \sum_{i=1}^N \|f(x_{a,i}) - f(x_{v,i})\|_2^2. \quad (3)$$

However, achieving seamless alignment requires more than just element-wise agreement. To address this, we impose a consistent alignment constraint. Consider normalized and standardized embeddings  $\mathbf{Z}_a, \mathbf{Z}_v \in \mathbb{R}^{T \times D}$ . Each embedding  $\mathbf{Z}_i$  is first normalized appropriately (mean subtraction, standard deviation division, and row-wise  $\ell_2$  normalization). We measure the difference between these correlation matrices using the Frobenius norm:

$$L_{cos} = \frac{1}{N} \sum_{i=1}^N \|\text{Corr}(\mathbf{Z}_{a,i}) - \text{Corr}(\mathbf{Z}_{v,i})\|_F. \quad (4)$$

Minimizing  $L_{cos}$  ensures that real and synthetic embeddings share not only similar feature magnitudes but also analogous internal semantic structures. Combining these two terms:

$$L_{msc} = \alpha L_{mse} + (1 - \alpha) L_{cos}. \quad (5)$$

**2) Representation Domain Alignment Loss:** To bridge the gap between real and synthetic distributions, a contrastive objective is employed to facilitate the learning of domain-invariant features. Specifically, we anchor on near-frontal (or minimally rotated) synthetic samples of the same speaker as positive samples, pairing them with corresponding real samples. Negative samples are constructed using embeddings from other data within the same batch. This setup encourages the model to focus on critical speech-related features, such as lip motion, while disregarding irrelevant factors like background, texture, and lighting. Formally, let  $f(\cdot)$  be the embedding function, and  $\{x_{a,i}, x_{v,i}\}_{i=1}^N$  be a real-synthetic pair of the same speaker with minimal head rotation differences. Positive similarity is defined as:

$$\text{pos\_sim} = \frac{\|f(x_{a,i}) - f(x_{v,i})\|_2}{\|f(x_{a,i})\|_2 + \|f(x_{v,i})\|_2} \quad (6)$$

For negatives, we use samples  $\{x_{a,j}\}_{j=1}^N$  from other data, treated as dissimilar:

$$\text{neg\_sim} = \frac{\|f(x_{a,i}) - f(x_{a,j})\|_2}{\|f(x_{a,i})\|_2 + \|f(x_{a,j})\|_2} \quad (7)$$

$L_{da}$  becomes:

$$L_{da} = -\frac{1}{N} \sum_{i=1}^N \frac{\exp(\text{pos\_sim})}{\exp(\text{pos\_sim}) + \exp(\text{neg\_sim})}. \quad (8)$$

where  $\tau$  is a temperature parameter and  $M$  is the number of negative samples. To align with downstream phoneme or voice categories for effective lip-reading, we integrate a Masked Multi-Modal Prediction Loss ( $L_{map}$ ). The MVL modeling training objective combines the three losses:

$$L_{MVL} = \lambda_{msc} L_{msc} + \lambda_{da} L_{da} + \lambda_{map} L_{map}, \quad (9)$$

where  $\lambda_{msc}, \lambda_{da}$ , and  $\lambda_{map}$  are weighting factors.

## C. Unified Modality Adapter

The UMA module consists of transformer layers, each receiving two inputs: the output of the corresponding layer of the MVL encoder (only audio feature) and the output of the previous UMA layer. This architectural setup ensures a progressive integration of the audio-visual features towards audio-visual-like representations. Let  $E_a$  denote the UMA-enhanced audio embeddings and  $E_v$  the corresponding audio-visual encoder embeddings for the same input context. We define a combined loss:

$$L_{ema} = \gamma L_{map} + (1 - \gamma) L_{eq}, \quad (10)$$

where  $L_{map}$  is the masked multi-modal prediction loss and  $L_{eq}$  is the feature similarity loss. We designate a set of early UMA layers as  $P$  and deeper layers as  $Q$ . For layers in  $P$ ,

we align the internal correlation structures between  $E_{a,i}$  and  $E_{v,i}$  to ensure semantic consistency, while for layers in  $Q$ , we directly minimize an element-wise feature discrepancy:

$$L_{eq} = \sum_{p \in P} \frac{1}{N} \sum_{i=1}^N \|\text{Corr}(E_{a,i}^{(p)}) - \text{Corr}(E_{v,i}^{(p)})\|_F^2 + \sum_{q \in Q} \frac{1}{N} \sum_{i=1}^N \|E_{a,i}^{(q)} - E_{v,i}^{(q)}\|_2^2. \quad (11)$$

where  $\lambda_p$  is a weighting factor controlling the importance of correlation alignment versus direct feature-level alignment.

## IV. EXPERIMENTS

### A. Experimental Settings

We evaluate our approach on LRSS3 [3] and OuluV52 [41]. Two subsets of LRSS3 are considered: a 30-hour subset (30h) and the full 433-hour set (433h). During 3D reconstruction, hyperparameters for geometry and texture refinement are tuned to preserve crucial lip motion features. Specifically, we use 150 epochs for geometry offset optimization, 80 for texture refinement, and 50 for joint geometry-texture optimization. Angles range from  $-25^\circ$  to  $25^\circ$ , sampled every  $5^\circ$ , with random angles in  $[-10^\circ, 10^\circ]$  for subtle variations. After synthesis, we form multi-view extended sets, ms30h: 30h set combining 30% synthetic and 70% real data, ms433h: 433h set combining 40% synthetic and 60% real data. OuluV52, which comprises more than 20k video recordings captured simultaneously from five different views ( $0^\circ, 30^\circ, 45^\circ, 60^\circ, 90^\circ$ ), is utilized to evaluate generalization under more realistic and challenging conditions. For modality-missing evaluation, we randomly mask the video modality at rates from 10% to 100%.

### B. Training Details

**Pretraining:** We initialize from the 5th iteration AV-HuBERT baseline, which provides trained pseudo-labels for the 6th iteration. This baseline is trained and fine-tuned following official protocols [15]. Our MVL model's encoder consists of 12 transformer blocks, and adopt a ResNet18-based front-end for visual feature extraction. Lip videos are resized to  $88 \times 88$  with a single channel. For the 30h set, during the first half of the update, we set  $\lambda_{msc} = 0.8, \lambda_{da} = 0.1, \lambda_{map} = 0.7$ . In the latter half, as the model gradually learns multi-view representations,  $\lambda_{msc}$  and  $\lambda_{da}$  are gradually reduced, stopping once the negative sample similarity in  $\lambda_{da}$  falls below 0.002. For the 433h data, the same schedule applies, but  $\lambda_{msc}$  and  $\lambda_{da}$  weights are reduced over the first third of the updates. When training the UMA module, the encoder weights remain frozen. Audio-only inputs pass through both the encoder and the UMA adapter. We set  $\gamma = 0.8$  and  $\lambda_p = 0.2$ .

**Fine-tuning:** For fine-tuning the MVL model with visual modality only, we use an attention-based sequence-to-sequence cross-entropy loss. For UMA model fine-tuning, we include both modalities, apply video masking at rates of  $\{0.0, 0.1, 0.2, \dots, 1.0\}$ , expanding 30 hours of original data into a total of 330 hours. All other hyperparameters remain consistent

TABLE 1. Comparison of our approach and AV-HuBERT. “0” corresponds to the standard LRSS3 test subset. For VSR, we incorporate multi-view synthetic data with angles of  $\pm 5^\circ$  or  $\pm 10^\circ$ . For AVSR, we apply 10% and 30% masking rate on video frames. † indicates that AV-HuBERT results are reproduced using the official code from [15], [24].

Model	Labelled data	Unlabelled data	Task	A	WER (%)	AI	AI	AI
	30	30	VSR	81.3	89.2	-	-	-
	30	433	VSR	84.3	69.0	72.1	-	-
	433	433	VSR	85.3	62.2	69.8	-	-
AV-HuBERT†	30	30	AVSR	5.7	-	-	6.8	6.7
	30	433	AVSR	5.8	-	-	6.1	6.1
Ours (MVL)	ms30	ms30	VSR	74.3	86.4	90.9	-	-
	ms30	ms433	VSR	75.9	68.8	74.8	-	-
	ms433	ms433	VSR	76.9	62.9	71.8	-	-
Ours (MVL+UMA)	ms30h	ms30	AVSR	5.6	-	-	5.2	5.7
	ms30h	ms433	AVSR	5.8	-	-	5.2	5.2

with the AV-HuBERT 5th iteration baseline to ensure a fair comparison.

### C. Experimental Results and Analysis

**1) Overall Performance:** Table 1 presents a comprehensive comparison of our models with the baseline AV-HuBERT. Without additional rotations or modality perturbations, the baseline AV-HuBERT achieves a WER of 5.8% when using 30 labelled hours and 433 unlabeled hours for AVSR tasks. By integrating the proposed MVL model, which leverages multi-view data and self-supervised losses ( $L_{msc}, L_{da}, L_{map}$ ), the WER further improves. Specifically, for the same data setting, our MVL model attains a lower WER of 43.2% for lipreading. Further incorporating the UMA module bolsters robustness against missing visual cues, maintaining a 5.6% WER in full audio-visual scenarios and showing stable performance even when video input is partially unavailable. In general, the first performance of the AVSR improves by about 16% relative to the baseline, highlighting the effectiveness of multi-view training and modality-invariant adaptation.

**2) Evaluation on Multi-View and Modality-Missing data:** We evaluate our model's performance under various multi-view conditions, including horizontal view angles up to  $\pm 30^\circ$ , combined yaw-pitch deviations, and partial modality-missing scenarios (show in Fig. 3). In the  $0^\circ$ – $20^\circ$  range, our approach achieves an overall 6% improvement in lip-reading performance, demonstrating superior robustness compared to the baseline. Even under more severe pose variations in both yaw and pitch, the MVL models maintain stable accuracy levels, consistently outperforming the baseline. When UMA is introduced to handle missing visual inputs, the model gains significant resilience to varying video inputs at 20% video masking rate, it surpasses the baseline by 16.7%. This gain stems not only from the robustness of UMA to incomplete modalities but also from the ability of MVL to extract richer lip-motion cues from diverse viewpoints, thereby enhancing overall audio-visual speech recognition capability.

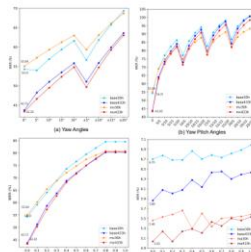


Fig. 3. Evaluation on multi-view and modality-missing data.

base30h and base433h denote the AV-HuBERT model finetuned on 30h and 433h datasets, ms30h and ms433h represent the MVL model finetuned on ms30h and ms433h. “Multi-View” indicates models pre-trained on ms433h data and fine-tuned on ms30 data, while “Mask†” signifies additional fine-tuning with masked samples.

**3) Generalization Performance Evaluation:** We conduct evaluations on four categories of genuine test data: three LRSS3 subsets grouped by increasing face yaw angles ( $<10^\circ$ ,  $10^\circ$ – $30^\circ$ ,  $>30^\circ$ ) and the OuluV52 dataset, which presents substantial viewpoint variations. As illustrated in Figure 4, when yaw angles stay within  $10^\circ$ , the MVL model trained on ms433h data achieves a WER of 31.1%, significantly outperforming the baseline. Even as yaw angles widen, our models continue to exhibit robust recognition performance on authentic multi-view samples. Notably, on OuluV52, which features profiles, large or even profile-angle views, the MVL approach leverages its strong multi-view lip-motion modeling to attain a WER of 43.25%, thereby showcasing its effectiveness under extreme pose conditions.

**4) Ablation Studies:** To further dissect the contributions of each component, we analyze performance under different configurations. Table II presents that without MVL, the baseline model struggles to align multi-view data, causing inconsistent gains or even degradations. Adding masking alone under a 30h fine-tuning setup provides limited improvements but does not fully capture multi-view representations. Notably, using only  $L_{msc}$  cannot effectively address background, illumination, or

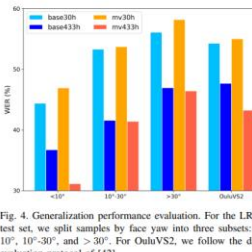


Fig. 4. Generalization performance evaluation. For the LRSS3 test set, we split samples by face yaw into three subsets:  $<10^\circ$ ,  $10^\circ$ – $30^\circ$ , and  $>30^\circ$ . For OuluV52, we follow the data evaluation protocol of [42].

TABLE II. Comparison among AV-HuBERT BASE models and our proposed models under diverse data and test settings on the test set. AVSR is pre-trained using 433h data and fine-tuned on 30h data. “Multi-View” indicates models pre-trained on ms433h data and fine-tuned on ms30 data, while “Mask†” signifies additional fine-tuning with masked samples.

Applies additional fine-tuning with masked samples.						
Method	Data		Loss			WER (%)
	Multi-View	Mask	$L_{MVC}$	$L_{rms}$	$L_{fs}$	
AV-HuBERT	✓	-	-	-	-	54.3
	✓	✓	✓	✓	✓	60.0
+ MVL	✓	-	✓	-	-	6.8
	✓	✓	✓	✓	✓	7.1
+ UMA	✓	-	✓	-	-	65.8
	✓	✓	✓	✓	✓	56.3
	✓	-	✓	✓	-	55.0
	✓	-	✓	-	✓	6.4
	✓	✓	✓	✓	✓	5.4

other extraneous factors, leading to suboptimal representation learning. Introducing  $L_{da}$  substantially enhances the model's multi-view capability (WER from 65.8% to 55.0%). Finally, equipping the model with UMA effectively mitigates the impact of missing or partially corrupted video signals. When multi-view data and the associated losses are collectively employed, the resulting model not only gains robustness to visual modality absence but also utilizes lip-motion information more effectively, yielding further improvements in audio-visual speech recognition.

## V. CONCLUSION

In this paper, we presented an enhanced self-supervised AVSR model that incorporates multi-view representation learning and modality-missing robustness. By leveraging 3D

facial avatar reconstruction, we generated multi-view speaker face videos, enabling the model to learn robust lip-reading capabilities across different head poses. The adapter module allows the model to handle scenarios where the visual modality is unreliable or missing by adjusting its reliance on visual features. Our experiments demonstrate that the proposed model consistently surpasses the AV-HuBERT model in both multi-view lip-reading and audio-visual speech recognition under missing-modality conditions. Furthermore, evaluations on real-world multi-view datasets validate the effectiveness of our approach in practical scenarios.

## REFERENCES

- [1] Eric David Perjan, “Automatic lipreading to enhance speech recognition (speech reading),” University of Illinois at Urbana-Champaign, 1984.
- [2] Benjamin Pétrotian, Peter Guez, and Eric Costas, “An end-to-end transformer approach for human based automatic lipreading,” in *Proc. ICIP*, IEEE, 1998, pp. 173–177.
- [3] Bowen Shi, Wei-Ning Han, Kohal Lakhotia, and Abdelrahman Mohamed, “Learning audio-visual speech representation by masked multi-modal cluster prediction,” in *Proc. ICML*, 2022.
- [4] Pan Zhen, Wenwen Yang, Wei Chen, Yanfeng Wang, and Jia Xu, “Learning from the mixture for end-to-end audio-visual speech recognition,” in *ICASSP*, IEEE, 2019, pp. 6500–6509.
- [5] Kenji Iwano, Tomoki Wakisaka, Satoshi Tamari, and Satoshi Furuta, “Audio-visual speech recognition using lip information from side-face images,” *EURASIP J. Audio Speech, and Music Process.*, vol. 2017, p. 2017, 2017.
- [6] Shihang Zhang, Ming Lei, Bin Ma, and Lei Xu, “Robust audio-visual speech recognition using bimodal fusion for automatic training and dropout regularization,” in *Proc. ICASSP*, IEEE, 2019, pp. 6570–6574.
- [7] Srikanth Parthasarathy and Shiva Sundaram, “Training strategies to handle missing modalities in audio-visual expression recognition,” in *Proc. ICMI*, 2020, pp. 400–403.
- [8] Philip William Grand, Mike Prindle, Tim Leister, Carsten Rohrer, Mathias Nollner, and Jochen Thies, “Neural head avatars from monocular video,” in *Proc. CVPR*, 2021, pp. 1851–1864.
- [9] Viktor Blanz and Thomas Vetter, “A morphable model for the synthesis of 3D faces,” *Proceedings of SIGGRAPH*, pp. 187–194, 1999.
- [10] Xianghui Chao, Xuebin Gong, Ming Cheng, Qi Deng, and Ming Li, “Cross-modal assisted training for abnormal event recognition in videos,” *Proceedings of the 2021 International Conference on Multimedia Interaction*, 2021.
- [11] Bowen Shi, Wei-Ning Han, Kohal Lakhotia, and Abdelrahman Mohamed, “Learning audio-visual speech representation by masked multi-modal cluster prediction,” *arXiv preprint arXiv:2301.02184*, 2023.
- [12] Christopher Bregler and Tachin Kung, “Gitterlip: for robust speech recognition,” in *Proc. ICASSP*, IEEE, 1994, pp. 11–60.
- [13] Joon Seon Chung, Andrew Senior, Oriol Vinyals, and Andrew Senior, “Lip reading sentences in the wild,” in *CVPR*, 2017, pp. 6447–6456.
- [14] Thomas Stylianidis and Georgios Tzimpanogiou, “Combining Residual Networks with LSTM for Lipreading,” in *Proc. Interspeech*, 2017, pp. 3652–3658.
- [15] Prithvankumar M. Suresh, Perisetti, and Maja Pantic, “Visual speech recognition for multiple languages in the wild,” *Nature Machine Intelligence*, vol. 4, no. 11, pp. 930–939, 2022.

- [16] Yiting Li, Yuki Takahashi, Tetsuya Takiguchi, and Yano Araki, “Lip reading using a dynamic feature of lip images and convolutional neural networks,” in *Proc. ICES*, IEEE, 2016, pp. 1–6.
- [17] Liqiang Nie, Mengzhou Xia, Xueming Song, Gonglu Wu, Hanyu Chang, and Jun Ge, “Multimodal activation: Advancing dialog robots without video words,” in *Proc. SIGDIAL*, pp. 491–500, 2021.
- [18] Prithvankumar M. Suresh, Perisetti, and Maja Pantic, “End-to-end audio-visual speech recognition with confounders,” in *Proc. ICASSP*, 2021, pp. 7613–7617.
- [19] Triantafyllos Alouzos, Joon Seon Chung, and Andrew Senior, “ASR is all you need: Cross-modal distillation for lip reading,” in *Proc. ICASSP*, 2022, pp. 2143–2147.
- [20] Bowen Shi, Wei-Ning Han, and Abdelrahman Mohamed, “Robust Self-Supervised Audio-Visual Speech Recognition,” in *Proc. Interspeech*, 2022, pp. 2118–2122.
- [21] Brandon Shillingford, Yannis Assaf, Matthew W. Hoffman, Thomas Pfister, et al., “Large-Scale Visual Speech Recognition,” in *Proc. Interspeech*, 2019, pp. 4135–4139.
- [22] Triantafyllos Alouzos, Joon Seon Chung, Andrew Senior, Oriol Vinyals, and Andrew Senior, “Deep audio-visual speech recognition,” in *Proc. CVPR*, 2018.
- [23] Mikaela Samarin, Stéphane Lathauwer, Sergey Tulyakov, Eliza Ricci, and Nica Sebe, “First order motion model for image animation,” in *Proc. NIPS*, 2019.
- [24] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tanzi, et al., “NeRF: Representing scenes as neural radiance fields for view synthesis,” in *Proc. ECCV*, 2020, pp. 405–421.
- [25] Zhenhua Li, Tao Yu, Hao Li, Guo-Xin Zhang, and Qinghui Di, “Learning a model of facial shape and expression from 4d scans,” in *Proc. SIGGRAPH Asia*, 2017.
- [26] Viktor Blanz and Thomas Vetter, “A morphable model for the synthesis of 3D faces,” in *Proc. SIGGRAPH*, 1999, pp. 187–194.
- [27] John Smith, Michael Zollhofer, Markus Stamminger, Christian Theobalt, and Mathias Nießner, “TheFaceReader: Real-time capture and reconstruction of 3d faces,” in *Proc. CVPR*, 2016, pp. 2387–2395.
- [28] Hongyong Chen, Pablo Garrido, Aydin Tavakoli, et al., “Deep video portraits,” *arXiv preprint arXiv:1805.1774*, 2018.
- [29] Seungwon Han, Yong Doo, Jaehoon G. Gwang, Han, and Seungwon Han, “Learning from the mixture for end-to-end audio-visual speech recognition for lip reading,” in *Proc. CVPR*, 2021, pp. 13325–13333.
- [30] Joon Seon Chung, Andrew Senior, Oriol Vinyals, and Andrew Senior, “Lip reading sentences in the wild,” *CVPR*, pp. 3444–3453, 2017.
- [31] Takaki Makino et al., “Recurrent neural network transfer for audio-visual speech recognition,” in