# UniSync: A Unified Framework for Audio-Visual Synchronization

Tao Feng[†,‡,*], Yifan Xie[†,*], Xun Guan[‡], Jiyuan Song[†], Zhou Liu[†], Fei Ma[†,**] and Fei Yu[†]

[†]Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen, China
[‡]Shenzhen International Graduate School, Tsinghua University, Shenzhen, China
Email: ft23@mails.tsinghua.edu.cn

*Abstract*—Precise audio-visual synchronization in speech videos is crucial for content quality and viewer comprehension. Existing methods have made significant strides in addressing this challenge through rule-based approaches and end-to-end learning techniques. However, these methods often rely on limited audio-visual representations and suboptimal learning strategies, potentially constraining their effectiveness in more complex scenarios. To address these limitations, we present UniSync, a novel approach for evaluating audio-visual synchronization using embedding similarities. UniSync offers broad compatibility with various audio representations (e.g., Mel spectrograms, HuBERT) and visual representations (e.g., RGB images, face parsing maps, facial landmarks, 3DMM), effectively handling their significant dimensional differences. We enhance the contrastive learning framework with a margin-based loss component and cross-speaker unsynchronized pairs, improving discriminative capabilities. UniSync outperforms existing methods on standard datasets and demonstrates versatility across diverse audio-visual representations. Its integration into talking face generation frameworks enhances synchronization quality in both natural and AI-generated content.

*Index Terms*—video analysis, audio-visual synchronization, representation learning, contrastive learning, talking face generation
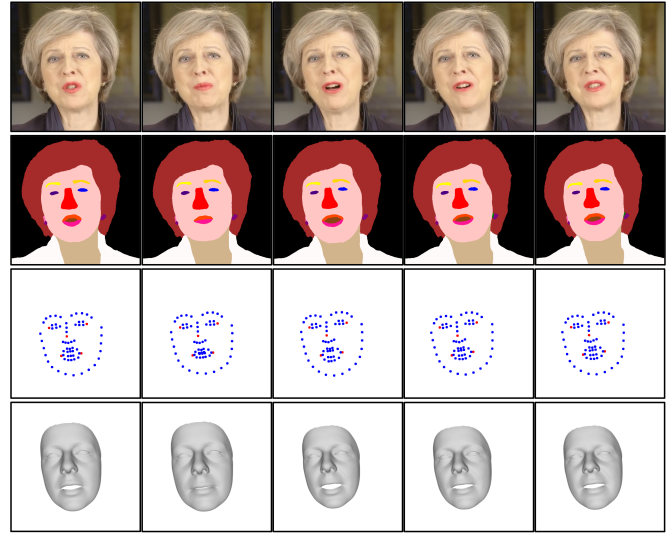
Fig. 1. Four most commonly used visual representation methods in speech video analysis. From top to bottom: RGB images, face parsing map, facial landmarks, and 3DMM.

## I. INTRODUCTION

Audio-visual synchronization, particularly lip synchronization, is critical in both real-world scenarios and AI-generated content (AIGC). Research [1] indicates that viewers can detect audio-visual desynchronization within a range of +45 to -125 milliseconds, highlighting the importance of precise synchronization. While various techniques like timecode encoding have been developed [2], methods that rely solely on audio and visual content for synchronization have emerged as more versatile solutions [3].

However, despite recent advancements in synchronization methods, existing approaches often rely on limited audio-visual representations, failing to leverage diverse representation techniques used in modern video analysis (as shown in Fig. 1). This limitation is particularly evident in talking face generation [4, 5], where various audio-visual representations have become essential for enhanced realism (summarized in Table I).

Furthermore, existing methods often employ inadequate negative sampling strategies [6–8], which can hinder their ability to effectively distinguish between synchronized and unsynchronized pairs [9]. This becomes particularly problematic when dealing with diverse real-world scenarios or AI-generated talking faces.

To address these challenges, we introduce UniSync, a novel method for evaluating audio-visual synchronization. The main contributions of this study are:

- We propose UniSync, a novel method compatible with multiple audio-visual representation techniques, and investigate the lip synchronization capabilities across different representation combinations.
- We enhance the model's performance by introducing a margin-based modification to the binary cross-entropy loss function and incorporating unsynchronized audio-visual pairs from different speakers as negative samples.
- We achieve state-of-the-art performance on both LRS2 [10] and CN-CVS [11] datasets and conduct compre-
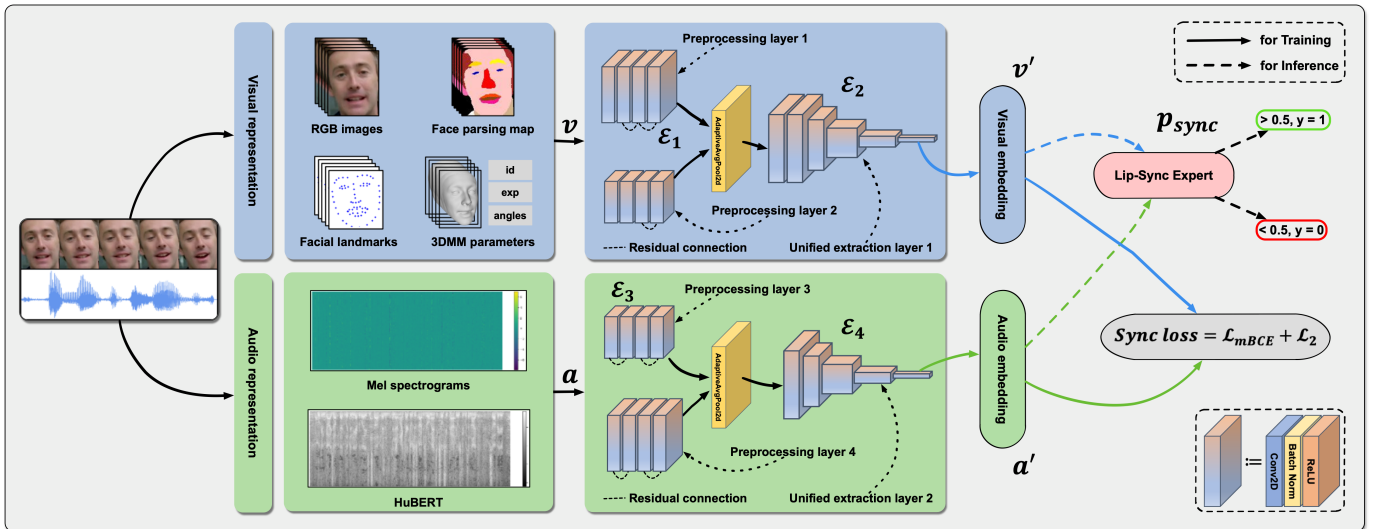
---

Fig. 2. The architecture of UniSync. For a given video segment, the process begins with the extraction of feature vectors using selected visual and audio representation methods, yielding $v$ and $a$ respectively. These vectors then undergo initial refinement through specialized preprocessing layers ($\varepsilon_1$ for visual data and $\varepsilon_3$ for audio data). Following an average pooling operation to standardize dimensions, unified feature extraction layers ($\varepsilon_2$ for visual data and $\varepsilon_4$ for audio data) generate the final embeddings $v'$ and $a'$. The degree of synchronization between visual and audio content is subsequently determined by computing the cosine similarity ($p_{\text{sync}}$) between these refined embeddings.

hensive experiments to demonstrate the versatility and effectiveness of our model.

## II. RELATED WORK

### A. Rule-Based Approaches

Rule-based approaches rely on predefined linguistic and phonetic rules to map speech sounds to visual representations. These methods typically involve a two-step process: first segmenting speech into phonemes, and then mapping these phonemes to visemes - the visual appearance of sounds on the lips [12].

Early works establish important foundations by using linear prediction techniques to decode phonemes from audio and correlate them with precise mouth positions [13]. Subsequent research advances this approach by categorizing English phonemes into 22 distinct visemes and implementing linear interpolation for smoother lip transitions [14]. Further developments propose that the relationship between phonemes and visemes is many-to-many due to visual coarticulation effects, leading to sophisticated mapping schemes using tree-based and K-means clustering techniques [15].

However, these methods often face limitations in real-world applications due to their reliance on controlled environments and extensive phonetic data [16–18]. These constraints have motivated the development of more flexible, learning-based approaches.

### B. End-to-End Learning Methods

End-to-end learning methods mark a significant advancement by directly learning from raw audio-visual data without relying on intermediate phonetic representations. Early progress emerges with HMM-based techniques that eliminate

TABLE I
DIVERSE REPRESENTATION TECHNIQUES USED IN RECENT TALKING FACE
GENERATION METHODS.

| Method | Year | Audio rep. | Visual rep. |
|--------|------|-----------|-------------|
| Wav2Lip [7] | 2020 | Mel spec. | RGB images |
| PC-AVS [24] | 2021 | Mel spec. | RGB images |
| AnyoneNet [25] | 2022 | Mel spec. | Landmarks |
| GeneFace [4] | 2023 | HuBERT | Landmarks, 3DMM |
| DreamTalk [26] | 2023 | HuBERT | RGB images, 3DMM |
| SegTalker [5] | 2024 | Mel spec. | Face parsing map |
| PointTalk [27] | 2024 | HuBERT | RGB images |

the need for explicit phoneme labeling [19]. The introduction of LSTM networks for direct audio feature processing [20] marks the beginning of deep learning's central role in this field.

A pivotal development comes with the introduction of dual-stream CNN architecture for feature extraction from Mel spectrograms and grayscale images [6]. This work inspires numerous subsequent developments, including enhanced models with multi-way matching and RGB image processing [8], cross-modal embedding matrices for accurate synchronization offset prediction [21], and multimodal transformer architectures for improved feature fusion [22].

However, these methods typically rely on limited types of audio-visual representations, primarily using grayscale or RGB images with Mel spectrograms. This constraint overlooks other valuable representation methods that have become crucial in modern video analysis, such as facial landmarks and 3DMM [23]. Our proposed UniSync addresses this limitation by supporting diverse visual and audio representations while maintaining strong performance across different frameworks.

## C. Applications in Talking Face Generation

In the field of talking face generation, lip synchronization plays a crucial role in enhancing content realism. Recent methods employ diverse audio-visual representations, as shown in Fig. 1. Notable examples include Wav2Lip [7] using Mel spectrograms with RGB images, and GeneFace [4] utilizing HuBERT features with facial landmarks. This diversity in representation methods, summarized in Table I, highlights the growing need for versatile synchronization approaches that can accommodate multiple input types.

## III. METHODOLOGY

### A. Unified Dual-Stream Architecture

UniSync employs a dual-stream architecture to process audio and visual representations separately. The visual stream supports four representation methods: RGB images, face parsing maps, facial landmarks, and 3DMM. These diverse representations capture different aspects of facial motion, from high-level appearance features to precise geometric information. The process involves a preprocessing layer ($\varepsilon_1$) for initial feature extraction, followed by an adaptive pooling layer to scale features to a fixed size. Larger-scale features like RGB images and face parsing maps undergo dimensionality reduction, while smaller-scale features are subjected to dimensionality increase during preprocessing. Finally, the scaled feature map is processed using a shared extraction layer ($\varepsilon_2$) to obtain the visual embedding ($v'$).

The audio stream processes both Mel spectrograms and HuBERT features through parallel pathways. Due to the inherent dimensional differences between these audio representations, we employ specialized preprocessing layers that adapt to their unique characteristics while preserving essential temporal information. Both streams ultimately produce identical-dimension embeddings, enabling direct computation of cosine similarity $p_{\text{sync}}$. Following established practices [7], a threshold of 0.5 serves as the decision boundary between synchronized and unsynchronized pairs.

UniSync's architecture balances simplicity and adaptability, using streamlined neural network components (2D convolutional layers, batch normalization, activation functions, and optional residual connections) to efficiently handle various audio-visual representation methods. The adaptive pooling layers play a crucial role by unifying inputs of different dimensions into a consistent format, enabling the subsequent shared extraction layers to process all representations uniformly. This architectural design maintains computational efficiency while supporting the flexibility needed for multi-modal synchronization tasks.

### B. Refined Learning Method for Audio-Visual Synchronization

We introduce a refined contrastive learning approach to enhance audio-visual synchronization detection. Our method combines an innovative loss function with cross-speaker negative samples [28], aiming to maximize cosine similarity for synchronized pairs while minimizing it for unsynchronized ones.
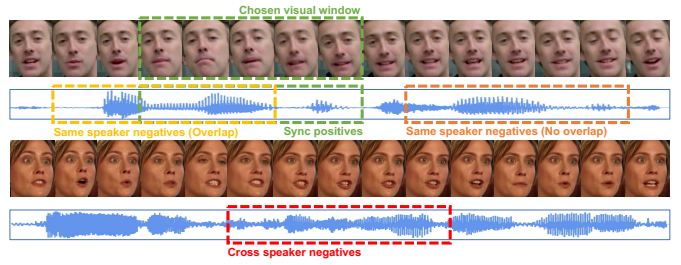


Fig. 3. Audio-visual contrastive learning samples: sync positives (matching audio and visual), same speaker negatives (with/without temporal overlap), and cross speaker negatives(mismatched speaker identity).

*1) Enhancing with Margin-Based Loss Function:* To improve the model's ability to discriminate between synchronized and unsynchronized audio-visual pairs, we develop an optimized loss function that combines a margin-based binary cross-entropy loss with $L_2$ regularization:

$$
\begin{aligned}
\mathcal{L}_{\text{ours}} = -\frac{1}{N} \sum_{i=1}^{N} \Big[ & y_i \log(p_{\text{sync},i}) + \\
& (1 - y_i) \log(1 - \max(0, \ p_{\text{sync},i} - m_i)) \Big] + \mathcal{L}_2
\end{aligned}
\tag{1}
$$

where $N$ is the total samples, $y_i$ the ground truth label, $p_{sync,i}$ the predicted synchronization probability, and $m_i$ the margin value. Traditional binary cross-entropy loss may not provide sufficient separation between synchronized and unsynchronized pairs. By introducing distinct margins for positive samples, same-speaker negatives, and cross-speaker negatives, we enforce a more rigorous decision boundary that helps the model better distinguish subtle differences in synchronization quality.

The synchronization probability $p_{\text{sync}}$ is computed using cosine similarity between audio embedding $a$ and visual embedding $v$, each representing a 0.2-second segment:

$$
p_{\text{sync}} = \cos(a, v) = \frac{a \cdot v}{\|a\|_2 \cdot \|v\|_2},
\tag{2}
$$

where $\cdot$ denotes the dot product of the vectors, and $\|a\|_2$ and $\|v\|_2$ are the Euclidean norms of the vectors.

Prior to computing the cosine similarity, both $a$ and $v$ undergo a ReLU activation function, ensuring that the calculated $p_{\text{sync}}$ values range between 0 and 1. This normalization process helps stabilize training and provides interpretable synchronization probabilities, where values closer to 1 indicate higher confidence in synchronization.

To mitigate overfitting and enhance generalization, we incorporate $L_2$ regularization:

$$
\mathcal{L}_2 = \lambda \|w\|_2^2 = \lambda \sum_{k=1}^{K} w_k^2
\tag{3}
$$

where $\lambda$ is a regularization parameter that controls the penalty strength, $K$ is the number of weights in the network, and $w_k$ represents each individual weight. This regularization term is particularly important given the diverse nature of our

| Year | Models | Backbone | Audio rep. | Acc.(%) |
|---|---|---|---|---|
| 2016 | SyncNet [6] | CNN | MFCC | 75.8 |
| 2018 | PM [8] | CNN | Mel spec. | 88.1 |
| 2020 | Wav2Lip [7] | CNN | Mel spec. | 90.7 |
| 2021 | AVST [29] | Transformer | Mel spec. | 92.0 |
| 2022 | VocaLiST [22] | Transformer | Mel spec. | **92.8** |
| 2024 | UniSync (Ours) | CNN | Mel spec. | 91.77 |
| 2024 | UniSync (Ours) | CNN | HuBERT | **94.27** |

| Model | FLOPs | Params | Speed* | Conv. Speed** | Acc. (%) |
|---|---|---|---|---|---|
| Wav2Lip [7] | 1.22G | 16.4M | 30s | Fast (30) | 90.7 |
| Vocalist [22] | 20.17G | 80.1M | 4.25min | Slow (289) | 92.8 |
| UniSync (Ours) | 1.73G | 16.3M | 33s | Fast (26) | **94.27** |

*Training time per epoch on RTX 4090 GPU
**Convergence speed to 80% val accuracy (epochs needed)

| Audio rep. | Visual rep. | Max Acc.(%) | Average Acc.(%) |
|---|---|---|---|
| HuBERT | RGB images | **93.25** | **90.09** |
| HuBERT | 3DMM | 91.40 | 89.11 |
| Mel spec. | RGB images | 91.31 | 87.23 |
| Mel spec. | 3DMM | 89.00 | 85.93 |
| HuBERT | Landmarks | 89.00 | 85.26 |
| HuBERT | Face parsing map | 86.04 | 82.50 |
| Mel spec. | Landmarks | 86.23 | 82.39 |
| Mel spec. | Face parsing map | 83.64 | 79.72 |

audio-visual representations and the need to maintain robust performance across different domains.

*2) Improving with Cross-Speaker Negative Samples:* To further enhance our proposed framework, we focus on effective sample selection, which plays a crucial role in model performance [9]. Traditional synchronization methods typically construct negative samples by shifting the audio or visual content within a single speaker's video. However, this approach may be insufficient for modern audio-visual tasks, particularly in talking face generation [4, 5], where the fundamental goal is to synchronize target audio with different speakers' facial movements.

Drawing inspiration from previous research [30], we distinguish between two types of negative samples: those from the same speaker and those from different speakers (see Fig. 3). By incorporating cross-speaker negative samples, where audio from one speaker is paired with visuals from another, we create a more challenging and realistic training scenario. The effectiveness of this strategy is particularly evident when applying UniSync to talking face generation tasks, where it significantly improves the synchronization quality of generated content (as demonstrated in Table VI). Our implementation maintains an appropriate proportion of cross-speaker negative samples in the training process, enabling the model to better handle the diverse synchronization scenarios encountered in real-world applications.

## IV. EXPERIMENTS

### A. Datasets and Implementation Details

To develop a unified framework compatible with various representation methods and generalizable across diverse speech videos, we utilize two large-scale lip-reading corpora: LRS2 [10] and CN-CVS [11]. LRS2, sourced from BBC broadcasts, contains over 224 hours of English spoken content from 1,000+ speakers. CN-CVS, compiled from various Chinese sources, comprises 273 hours of Mandarin visual-speech from 2,500+ speakers. Both datasets offer diverse content for lip-sync tasks, encompassing different speaking styles, emotional expressions, and recording environments.

Our UniSync model is trained using a batch size of 64 on an NVIDIA RTX 4090 GPU with 24GB of memory. The training process employs the Adam optimizer [31] with a learning rate of $1 \times 10^{-4}$ and momentum parameters $\beta_1 = 0.5$, $\beta_2 = 0.999$. To mitigate overfitting, we incorporate weight decay with a coefficient of $1 \times 10^{-4}$.

### B. Baselines and Evaluation Metrics

For comprehensive evaluation, we compare our UniSync against several methods: SyncNet [6], which pioneers the CNN-based approach with MFCC features; Perfect Match (PM) [8] and Wav2Lip [7], which advance the field using Mel spectrograms; and more recent Transformer-based models like AVST [29] and VocaLiST [22].

We employ two primary evaluation metrics: Lip-Sync Accuracy and Lip Sync Error, following the methodology proposed by [7]. Lip-Sync Accuracy is calculated using randomly selected synchronized and unsynchronized audio-visual pairs. Additionally, we utilize LSE-D (Lip Sync Error - Distance) and LSE-C (Lip Sync Error - Confidence), based on the pretrained SyncNet model [6], to evaluate synchronization accuracy in uncontrolled environments. These industry-standard metrics enable direct performance comparisons with existing and future methods.

### C. Comparison of Lip-Sync Accuracy with Popular Models

As shown in Table II, experimental results demonstrate the effectiveness of the proposed model, particularly when using HuBERT as audio representation. Our CNN-based UniSync achieves an impressive accuracy of 94.27%, surpassing the performance of recent transformer-based models such as AVST [29] and VocaLiST [22]. This result underscores the potential of combining traditional CNN architectures with state-of-the-art audio processing techniques. To ensure a fair comparison with previous methods, we also evaluate UniSync using Mel spectrograms as the audio input. In this configuration, UniSync achieves an accuracy of 91.77%, which, while slightly lower than the transformer-based methods, still represents a significant improvement over earlier CNN-based approaches like PM [8] and Wav2Lip [7].

| Aspect | Value | Acc.(Mel) | Acc.(HuBERT) |
|---|---|---|---|
| Cross-speaker proportion[1] | 0.0 | 88.96 | 91.79 |
| | **0.2** | **89.23** | **92.18** |
| | 0.5 | 88.93 | 91.96 |
| Margin (same, cross)[2] | (0, 0) | 91.22 | 93.62 |
| | (0.1, 0.3) | **91.77** | 93.99 |
| | (0.3, 0.7) | 91.50 | **94.27** |

[1] Results reported as mean accuracy across multiple well-converged runs, ensuring a reliable measure of overall performance.
[2] Results indicate maximum accuracy achieved, demonstrating the optimal potential of each configuration.

| Method | LRS2 | | CN-CVS | | GeneFace's dataset[1] | |
|---|---|---|---|---|---|---|
| | LSE-D ↓ | LSE-C ↑ | LSE-D ↓ | LSE-C ↑ | LSE-D ↓ | LSE-C ↑ |
| LipGAN [32] | 10.330 | 3.199 | - | | - | |
| Wav2Lip [7] | 7.521 | 6.406 | 11.025 | 1.439 | - | |
| **Wav2Lip + UniSync** | **6.647** | **6.842** | **9.898** | **2.2512** | - | |
| GeneFace [4] | - | | - | | 9.809 | 4.765 |
| **GeneFace + UniSync** | - | | - | | **9.396** | **5.294** |

[1] GeneFace's dataset refers to a video of a female speaker approximately 4 minutes long, some frames of which are shown in Fig. 1.

Furthermore, as illustrated in Table III, UniSync demonstrates superior computational efficiency compared to transformer-based alternatives. Our model maintains comparable computational complexity to Wav2Lip [7] while significantly outperforming VocaLiST [22] in terms of both parameter count and FLOPs. The training efficiency is especially significant, with UniSync requiring substantially less time per epoch and fewer epochs to achieve convergence compared to VocaLiST [22].

These results demonstrate the effectiveness of our approach. To further understand the impact of different representation choices, we conduct a comprehensive analysis of various audio-visual combinations in the following section.

### D. Evaluating Lip-Sync Accuracy Across Audio-Visual Representations

We evaluate the performance of various audio-visual representation combinations, as shown in Table IV. The results reveal several important patterns: (1) HuBERT consistently outperforms Mel spectrograms across all visual representations, with improvements ranging from 2.4% to 3.5% in average accuracy; (2) Among visual representations, RGB images achieve the best performance (93.25%), followed by 3DMM (91.40%), while face parsing maps show relatively lower accuracy; (3) The combination of HuBERT and RGB images yields both the highest maximum accuracy (93.25%) and average accuracy (90.09%), suggesting this pairing as the optimal choice for practical applications. Notably, the performance gap between different visual representations is more pronounced when using Mel spectrograms compared to HuBERT, indicating HuBERT's stronger capability in extracting discriminative audio features.

### E. Ablation Study

We conduct ablation studies on two critical components of UniSync: the proportion of cross-speaker negative samples and the margin values in the loss function. Table V presents our findings using both Mel spectrogram and HuBERT audio representations.

For cross-speaker negative samples, we observe that a moderate proportion of 0.2 achieves optimal performance, improving accuracy from 88.96% to 89.23% with Mel spectrograms and from 91.79% to 92.18% with HuBERT. Higher proportions (0.5) show diminishing returns, suggesting that a balanced mix of same-speaker and cross-speaker negatives is crucial for effective training.

The introduction of margin values in the loss function consistently improves performance compared to the baseline without margins. As shown in Table V, both audio representations benefit from this enhancement, with HuBERT achieving slightly better results overall.

### F. Improving Lip-Sync Quality in Talking Face Generation

We integrated UniSync into two representative talking face generation models: Wav2Lip [7] and GeneFace [4]. As shown in Table VI, UniSync significantly improves lip synchronization quality across different scenarios. For Wav2Lip, our integration reduces LSE-D from 7.521 to 6.647 on LRS2 (English) and from 11.025 to 9.898 on CN-CVS (Mandarin), while simultaneously increasing LSE-C scores. Similar improvements are observed with GeneFace, where LSE-D decreases from 9.809 to 9.396 and LSE-C increases from 4.765 to 5.294 on real-world video content.

These consistent improvements across different models, languages, and datasets demonstrate UniSync's effectiveness as a general-purpose lip synchronization enhancement module. Notably, the performance gains in both controlled (LRS2) and real-world scenarios (GeneFace's dataset) suggest UniSync's strong adaptability to various practical applications.

## V. CONCLUSION

We present UniSync, a novel method for audio-visual synchronization that addresses key limitations of existing approaches. UniSync's broad compatibility with various representation techniques, including RGB images, face parsing maps, facial landmarks, and 3DMM for visual content, alongside Mel spectrograms and HuBERT features for audio, enhances its adaptability. Our refined contrastive learning strategy, incorporating both same-speaker and cross-speaker negative samples, substantially improves synchronization evaluation capability. Through extensive experiments, UniSync demonstrates superior performance in standard benchmarks and enhances lip synchronization quality when integrated into talking face generation models. The method's ability to handle diverse representations while maintaining computational efficiency makes it valuable for real-world applications and future audio-visual tasks, especially in the field of AI-generated content.

REFERENCES

[1] ITU Radiocommunication, "Relative timing of sound and vision for broadcasting," .

[2] T. Amyes, "Film in audio post production," in *Audio Post Production for Television and Film*, pp. 68–76. Routledge, 2013.

[3] X. Li, X. Wang, K. Wang, and S. Lian, "A novel speech-driven lip-sync model with cnn and lstm," in *2021 14th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*. IEEE, 2021, pp. 1–6.

[4] Z. Ye, Z. Jiang, Y. Ren, J. Liu, J. He, and Z. Zhao, "Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis," *arXiv preprint arXiv:2301.13430*, 2023.

[5] L. Xiong, X. Cheng, J. Tan, X. Wu, X. Li, L. Zhu, F. Ma, M. Li, H. Xu, and Z. Hu, "Segtalker: Segmentation-based talking face generation with mask-guided local editing," *arXiv preprint arXiv:2409.03605*, 2024.

[6] J. S. Chung and A. Zisserman, "Out of time: automated lip sync in the wild," in *Computer Vision–ACCV 2016 Workshops: ACCV 2016 International Workshops, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part II 13*. Springer, 2017, pp. 251–263.

[7] K. R. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, and C. V. Jawahar, "A lip sync expert is all you need for speech to lip generation in the wild," in *Proceedings of the 28th ACM international conference on multimedia*, 2020, pp. 484–492.

[8] S. W. Chung, J. S. Chung, and H. G. Kang, "Perfect match: Improved cross-modal embeddings for audio-visual synchronisation," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3965–3969.

[9] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus, "Hard negative mixing for contrastive learning," *Advances in neural information processing systems*, vol. 33, pp. 21798–21809, 2020.

[10] T. Afouras, J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, "Deep audio-visual speech recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 12, pp. 8717–8727, 2018.

[11] C. Chen, D. Wang, T. F. Zheng, and T. Fang, "Cn-cvs: A mandarin audio-visual dataset for large vocabulary continuous visual to speech synthesis," in *2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[12] H. L. Bear and R. Harvey, "Phoneme-to-viseme mappings: the good, the bad, and the ugly," *Speech Communication*, vol. 95, pp. 40–67, 2017.

[13] J. Lewis, "Automated lip-sync: Background and techniques," *The Journal of Visualization and Computer Animation*, vol. 2, no. 4, pp. 118–122, 1991.

[14] S. Morishima, S. Ogata, K. Murai, and S. Nakamura, "Audio-visual speech translation with automatic lip sync-qronization and face tracking based on 3-d head model," in *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing*. IEEE, 2002, vol. 2, pp. II–2117.

[15] W. Mattheyses, L. Latacz, and W. Verhelst, "Comprehensive many-to-many phoneme-to-viseme mapping and its application for concatenative visual speech synthesis," *Speech Communication*, vol. 55, no. 7-8, pp. 857–876, 2013.

[16] L. Wang, X. Qian, W. Han, and F. K. Soong, "Synthesizing photo-real talking head via trajectory-guided sample selection," in *Eleventh Annual Conference of the International Speech Communication Association*, 2010.

[17] R. Anderson, B. Stenger, V. Wan, and R. Cipolla, "An expressive text-driven 3d talking head," in *ACM SIGGRAPH 2013 Posters*, pp. 1–1. 2013.

[18] B. Fan, L. Wang, F. K. Soong, and L. Xie, "Photo-real talking head with deep bidirectional lstm," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 4884–4888.

[19] M. Brand, "Voice puppetry," in *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999, pp. 21–28.

[20] T. Shimba, R. Sakurai, H. Yamazoe, and J. H. Lee, "Talking heads synthesis from audio with deep neural networks," in *2015 IEEE/SICE International Symposium on System Integration (SII)*. IEEE, 2015, pp. 100–105.

[21] Y. J. Kim, H. S. Heo, S. W. Chung, and B. J. Lee, "End-to-end lip synchronisation based on pattern classification," in *2021 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2021, pp. 598–605.

[22] V. S. Kadandale, J. F. Montesinos, and G. Haro, "Vocalist: An audio-visual synchronisation model for lips and voices," *arXiv preprint arXiv:2204.02090*, 2022.

[23] L. Chen, G. Cui, Z. Kou, H. Zheng, and C. Xu, "What comprises a good talking-head video generation?: A survey and benchmark," *arXiv preprint arXiv:2005.03201*, 2020.

[24] H. Zhou, Y. Sun, W. Wu, C. C. Loy, X. Wang, and Z. Liu, "Pose-controllable talking face generation by implicitly modularized audio-visual representation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 4176–4186.

[25] X. Wang, Q. Xie, J. Zhu, L. Xie, and O. Scharenborg, "Anyonenet: Synchronized speech and talking head generation for arbitrary persons," *IEEE Transactions on Multimedia*, vol. 25, pp. 6717–6728, 2022.

[26] Y. Ma, S. Zhang, J. Wang, X. Wang, Y. Zhang, and Z. Deng, "Dreamtalk: When expressive talking head generation meets diffusion probabilistic models," *arXiv preprint arXiv:2312.09767*, 2023.

[27] Y. Xie, T. Feng, X. Zhang, X. Luo, Z. Guo, W. Yu, H. Chang, F. Ma, and F. R. Yu, "Pointtalk: Audio-driven dynamic lip point cloud for 3d gaussian-based talking head synthesis," *arXiv preprint arXiv:2412.08504*, 2024.

[28] A. van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[29] H. Chen, W. Xie, T. Afouras, A. Nagrani, A. Vedaldi, and A. Zisserman, "Audio-visual synchronisation in the wild," *arXiv preprint arXiv:2112.04432*, 2021.

[30] B. Korbar, D. Tran, and L. Torresani, "Cooperative learning of audio and video models from self-supervised synchronization," *Advances in Neural Information Processing Systems*, vol. 31, 2018.

[31] D. P. Kingma, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[32] K. R. Prajwal, R. Mukhopadhyay, J. Philip, A. Jha, V. Namboodiri, and C. V. Jawahar, "Towards automatic face-to-face translation," in *Proceedings of the 27th ACM international conference on multimedia*, 2019, pp. 1428–1436.