# A Survey on Deep Multi-modal Learning for Body Language Recognition and Generation

Li Liu, *Member, IEEE*, Lufei Gao, Wentao Lei, Fengji Ma, Xiaotian Lin, Jinting Wang

**Abstract**—Body language (BL) refers to the non-verbal communication expressed through physical movements, gestures, facial expressions, and postures. It is a form of communication that conveys information, emotions, attitudes, and intentions without the use of spoken or written words. It plays a crucial role in interpersonal interactions and can complement or even override verbal communication. Deep multi-modal learning techniques have shown promise in understanding and analyzing these diverse aspects of BL, which often incorporate multiple modalities. The survey explores recent advances in deep multi-modal learning, emphasizing their applications to BL generation and recognition. Several common BLs are considered *i.e.*, Sign Language (SL), Cued Speech (CS), Co-speech (CoS), and Talking Head (TH), and we have conducted an analysis and established the connections among these four BL for the first time. Their generation and recognition often involve multi-modal approaches, for example, multi-modal feature representation, multi-modal fusion, and multi-modal joint learning will be introduced. Benchmark datasets for BL research are well collected and organized, along with the evaluation of state-of-the-art (SOTA) methods on these datasets. The survey highlights challenges such as limited labeled data, multi-modal learning, and the need for domain adaptation to generalize models to unseen speakers or languages. Future research directions are presented, including exploring self-supervised learning techniques, integrating contextual information from other modalities, and exploiting large-scale pre-trained multi-modal models. Real-world applications and user-centric evaluations are emphasized to drive practical adoption. In summary, this survey paper provides a comprehensive understanding of deep multi-modal learning for various BL generations and recognitions for the first time. By analyzing advancements, challenges, and future directions, it serves as a valuable resource for researchers and practitioners in advancing this field. In addition, we maintain a continuously updated paper list for deep multi-modal learning for BL recognition and generation: https://github.com/wentaoL86/awesome-body-language.

**Index Terms**—Deep Multi-modal Learning, Body Language, Sign Language, Cued Speech, Co-speech, Talking Head, Recognition and Generation.

---

## 1 INTRODUCTION

BODY language (BL), as a vital component of non-verbal communication, holds great significance in facilitating effective communication and enhancing social interactions. The ability to analyze and understand BL has various applications, ranging from BL recognition and generation to digital human interaction and assistive technologies. Understanding BL often necessitates the incorporation of multiple modalities. Deep multi-modal learning, which combines visual, audio and text modalities have emerged as a promising approach to enhancing the accuracy and robustness of intelligent BL multi-modal conversion systems.

In this survey, we primarily focus on four typical BLs and use them as examples to review and analyze the multi-modal BL recognition and generation. Figure 1 presents a simple diagram for the four types of BLs, *i.e.*, Cued Speech (CS) [1], Sign Language (SL) [2], Co-speech (CoS) [3] and Talking Head (TH) [4]. In this field, there have been numerous previous works, which have made significant progress. However, despite the progress made in deep multi-modal learning for BL generation and recognition, several challenges and open research questions remain, such as the multi-modal learning of different types of data modalities, the scarcity of labeled datasets, representing fine-grained cues, modeling temporal dynamics, and limited computational resources. These challenges need to be addressed in multi-modal BL recognition and generation to further advance
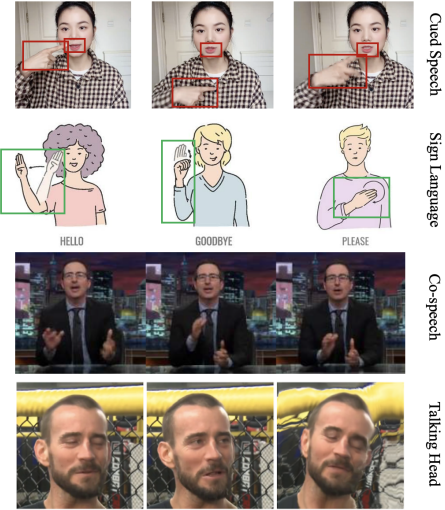


Fig. 1: Examples of Cued Speech, Sign Language, Co-speech and Talking Head, respectively.

the field and make applications in human-computer interaction (HCI), social robotics, and affective computing more effective, etc.

**Organization of This Survey.** In this survey, we first introduce four typical variants of BL and establish the connections between these four types in Section 2. Then, We organize and present various types of datasets for BL recognition and generation, along with evaluation metrics in Section 3. In Sections 4

---

- *Li Liu, Wentao Lei, Fengji Ma, Xiaotian Lin, and Jinting Wang are with the Hong Kong University of Science and Technology (Guangzhou), Guangzhou 511458, China. E-mail: avrillliu@hkust-gz.edu.cn.*
- *Lufei Gao is with the Shenzhen Research Institute of Big Data, Shenzhen, China.*

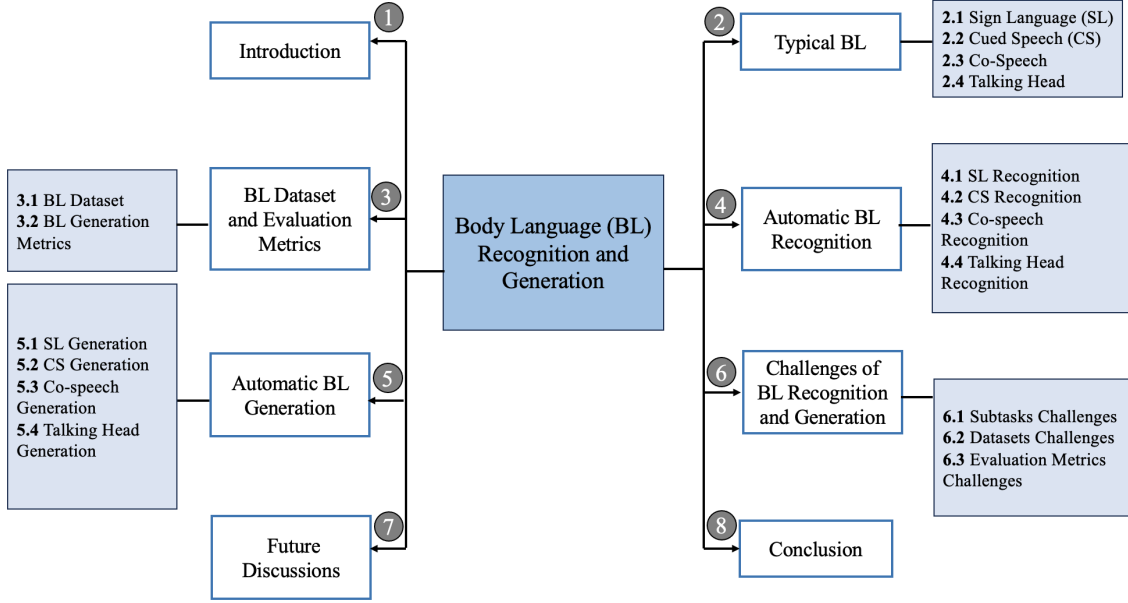*All authors are equal contribution.*

Fig. 2: The architecture of this survey.

and 5, we provide detailed reviews of the BL recognition and generation of CS, SL, CoS and TH, respectively. Furthermore, in Section 6, we give a detailed analysis of the challenges for these types of BL. Finally, we discuss and conclude this survey by proposing multiple research directions that need to be studied. The architecture of this survey is visualized in Figure 2. The structured taxonomy of the existing BL research and some representative works are shown in Figure 3.

TABLE 1: The number of existing reviews.

| Type | SL | CS | CoS | TH | LR | SL+CS | LR+TH | Total |
|------|----|----|-----|----|----|-------|-------|-------|
| R | 5 | 1 | 0 | 0 | 5 | 1 | 0 | 12 |
| G | 4 | 0 | 1 | 3 | 0 | 0 | 0 | 8 |
| R&G | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 2 |
| Total | 10 | 1 | 1 | 3 | 5 | 1 | 1 | 22 |

*The corresponding terms for the abbreviations are as follows: R – Recognition; G – Generation; SL – Sign Language; CS – Cued Speech; CoS – Co-speech; TH – Talking Head; LR – Lip Reading.*

**Differences from Existing Reviews.** Table 1 presents the number of review articles related to BL recognition and generation in the relevant field. While there are already 22 existing surveys, the differences between our survey and these prior works can be summarized as follows:

- **Scope**. Existing reviews on BL [5], [6], [7] have only focused on specific subtasks within the field. For BL recognition, the reviews [5], [8], [9], [10], [11], [12], [13], [14], [15], [16] concentrate on SL recognition. Regarding BL generation, [6] only explores CoS generation and [17], [18], [19] delves into TH generation. [20] integrates subtasks: LR recognition and TH generation. Unlike the reviews mentioned earlier, this paper focuses on two primary tasks: **recognition** and **generation**. Each task is expanded to incorporate four different types of BL: **SL**, **CS**, **CoS**, and **TH**. As far as we know, this is the

first to encompass all four types of BL along with their corresponding recognition and generation tasks.

- **Timeline.** This survey highlights the latest advances, major challenges and deep learning (DL)-based multi-modal approaches in the aforementioned research areas from 2017 to the present. Please note that we will consistently update the repository we maintain with the latest developments. It is expected that this study will facilitate knowledge accumulation and the creation of deep multi-modal BL methods, providing readers, researchers, and practitioners with a roadmap to guide future direction.

To summarize, this survey provides a thorough examination of the progress made in deep multi-modal learning techniques for automatic BL recognition and generation. It also outlines the road ahead for future research in this area. The objective is to offer researchers and practitioners a consolidated understanding of the field, covering the foundational principles, multi-modal fusion methods, DL architectures, benchmark datasets, challenges, and potential directions.

## 2 TYPICAL BODY LANGUAGE

BL through which humans convey information usually involves five aspects, *i.e.*, gestures, facial expressions, lip reading (LR), head pose and postures. In this survey, we refer to these five aspects as the basic elements of BL.

Gestures refer to the use of hand movements to convey meaning. People communicate through actions such as waving, pointing, or gesturing with their hands. Additionally, facial expressions play a crucial role as a basic element of BL. Humans express emotions and intentions by altering the facial muscles around the eyes, eyebrows, mouth, etc. Another fundamental element is LR, which involves interpreting speech by observing the movements of the lips and mouth. Furthermore, head pose, including tilting or turning the head, can also convey information related to attention, interest, or specific desires. Lastly,
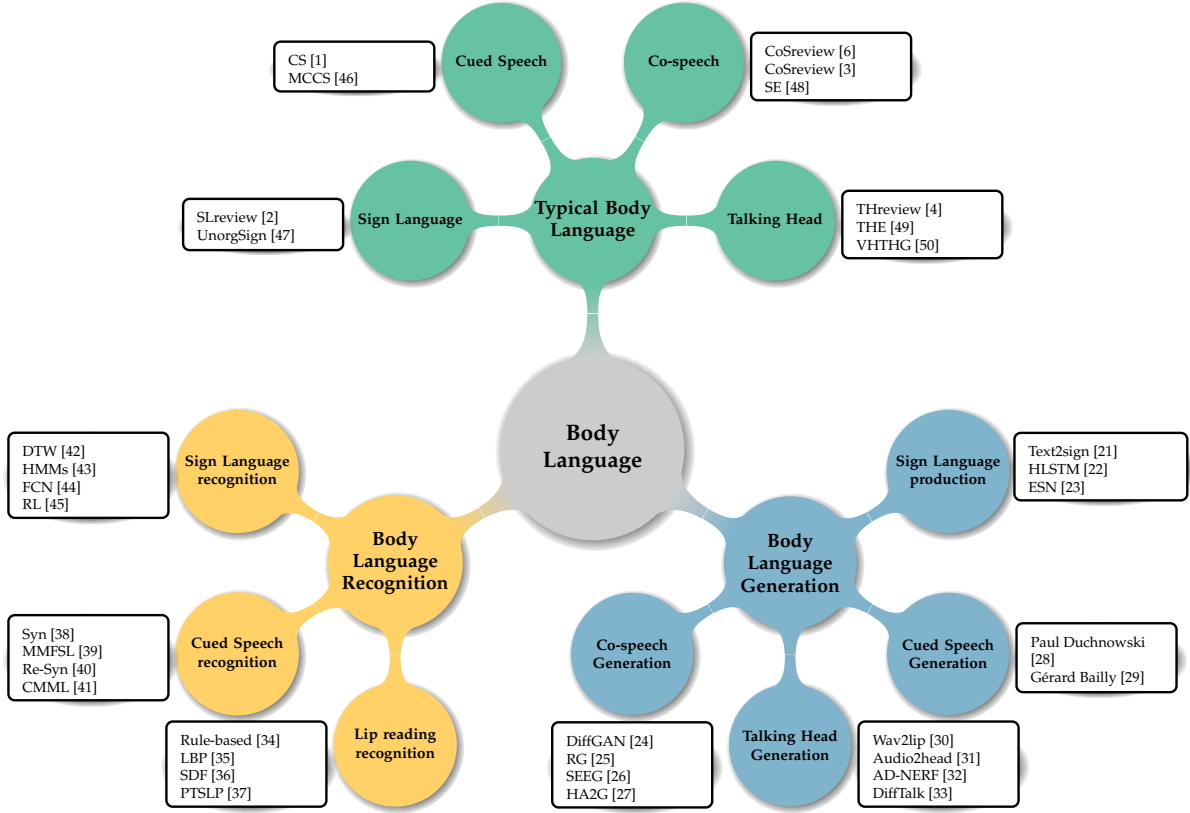
Fig. 3: Structured taxonomy of the existing BL research which includes three genres. Only several representative methods of each category are demonstrated.
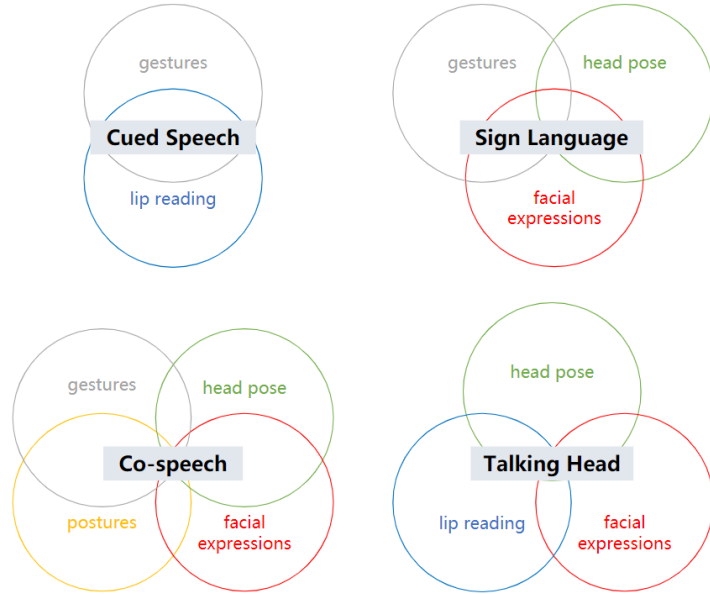


Fig. 4: Element compositions of four typical body language cases.

postures, such as standing, sitting, or body leaning, contribute to conveying emotional states and social intentions within BL.

It is common for BL cases to consist of two or more of these modalities. As shown in Figure 4, we listed four typical BL cases that are discussed in this survey, and each of them can be regarded as a composition of the basic BL elements.

In this section, we will provide a comprehensive overview of these four BL cases, including their concepts, significance, and the challenges that exist in their corresponding recognition or generation tasks.

## 2.1 Sign Language

SL is categorized as a natural language commonly used in deaf communities [2]. Based on data from the World Federation of the Deaf, the worldwide population of the deaf is estimated to be around 72 million, with over 80% living in developing nations [47]. Over 300 different SLs are used by these individuals, each having its own distinct vocabulary and grammar. SL is also known as a visual language which is generally composed of several visual partials, such as gestures, facial expressions, head pose and body postures. Specifically, six basic parameters are listed as the basic components of SL in [51], *i.e.*, hand shape, orientation, movement, location, mouth shape and eyebrow movements. Taking an overall perspective into account, we regard gestures, facial expressions, and head poses as the primary visual modalities in SL.

SL is the major communication tool for the deaf, yet it is difficult to be mastered. In order to eliminate communication barriers, it is of great significance to develop technologies for automatic SL processing, including SL recognition (SLR) that extracts words or utterances by capturing and analyzing image or video sequences of the SL data, SL generation (SLG) that generates visualizable SL animations from input with semantic meaning and SL translation that translates the extracted information to another signed or spoken language [52] [53]. This survey

mainly focuses on the literature review of SLR and SLG in order to deeply understand the important issues and difficulties in the field of SL processing.

As a highly dynamic and multi-modal visual language, SL involves a combination of multiple visual elements that have complementary semantics. Therefore, extracting and fusing high-dimensional features from different modalities effectively is an important task. Deep multi-modal learning techniques play a pivotal role in addressing these challenges and advancing the field of SL processing. By combining visual and spatial information from video or depth sensors with linguistic cues, these approaches have shown promising results in improving the accuracy and naturalness of SLR and SLG systems.

In Section 4.1 and 5.1, we investigate the recent research advancements and techniques in deep multi-modal learning specifically for SLR and SLG. Besides, we delve into the challenges that are associated with these tasks and emphasize the potential applications and future directions in this field.

## 2.2 Cued Speech

CS is a visual communication system proposed by Cornett [1] to enhance speech perception for individuals with hearing loss. CS uses a set of hand shapes and positions, named cues, to code the phonemes such as consonants and vowels. Figure 5 presents the chart for Mandarin Chinese CS (MCCS) [46], [54]. Gestures in CS functions as a complement to lip-reading, visualizing the phonetic details that can be observed from the mouth movements to remove ambiguities caused by lip-reading alone. As a clear and unambiguous visual counterpart to the auditory information in spoken language, CS enables individuals with hearing loss to better understand and distinguish speech sounds, facilitating their language acquisition, spoken capabilities, reading skills, and overall communication abilities.
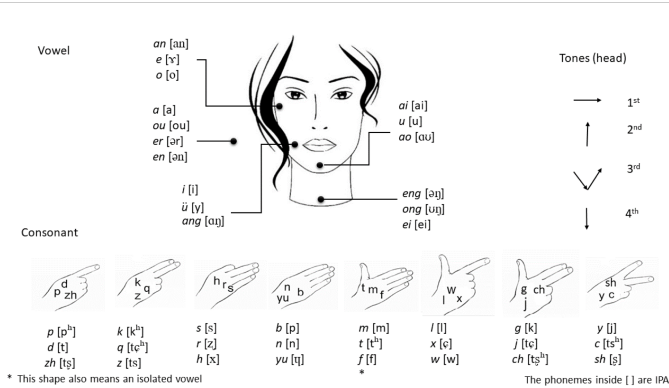


Fig. 5: The corresponding table between hand shapes and hand positions in Mandarin Chinese CS for vowels and consonants, respectively.

CS has currently been adapted to approximately 65 languages and dialects globally, including English, French, Chinese, etc. [55]. Recently, there has been growing interest in developing technologies for automatic recognition and generation in the CS research field [56], [57]. These technologies aim to enhance accessibility for people who primarily use CS for communication. For instance, utilizing automatic CS recognition (ACSR), people can effortlessly transcribe gestures and lip-reading into corresponding spoken language at a phonemic level [58], [59]. In the opposite direction, a digital agent equipped with automatic CS generation (ACSG) can convert spoken input into authentic CS expressions.

To process CS data efficiently, it is crucial to effectively extract information from two modalities: hand and lip movement. However, this task poses challenges in several aspects. Firstly, there exists an inherent asynchrony phenomenon when the human brain processes speech with gestures [40]; secondly, recognizing or generating appropriate hand shapes and lip movements entails tackling fine-grained image processing problems [60], [61], [62]. Consequently, deep multi-modal learning techniques have emerged as a prominent research trend to uncover the interplay between gestures and LR, aiming to achieve high-performance ACSR or ACSG systems.

In Section 4.2 and 5.2, we will discuss some SOTA methods to solve the problems lying in CS processing. We explore the challenges and opportunities in leveraging deep multi-modal learning techniques in this new research area, aiming to enhance the accessibility and inclusivity of communication for individuals who rely on CS.

## 2.3 Co-speech

CoS refers to the non-verbal behaviors and signals that accompany and complement spoken language during communication [6]. It encompasses various visual cues, such as gestures, body postures, and facial expressions such as eye gaze and blinking, which can be used in conjunction with speech to convey additional information and meaning [63] [64]. CoS gestures contribute substantially to the overall comprehension and interpretation of spoken language [65]. They serve as contextual cues, accentuate salient points, convey emotional states, and facilitate social interactions [66].

With the development of AI agents technologies, there has been extensive research exploration in CoS generation or synthesis to give AI agents such as digital humans more expressive and realistic BL [27], [67], [68], [69], [70]. The primary objective of this task is to generate a sequence of human BL by utilizing speech audio and transcripts as input, enhancing the performance of human-machine interaction systems. On the other hand, most existing gesture recognition methods primarily focus on recognizing specific types of gestures [71], [72], [73], [74], overlooking their connections with other modalities such as speech.

CoS signals not only play a crucial role in enhancing the clarity, expressiveness and emotional content of verbal communication, but also capture the rich communicative context, and reveal the speaker's social identity and cultural affiliation [48]. Therefore, it is a growing trend towards exploring multi-modal approaches that take into account both the visual information from gestures and the accompanying speech signals, which allows for more comprehensive and accurate analysis in the areas such as emotion recognition and dialogue understanding [75].

In Section 4.3, we provide a brief overview of the work on automatic CoS recognition. Due to its limited application scenarios, research in this field is relatively scarce. In Section 5.3, we review the SOTA techniques and advancements in deep multi-modal learning for CoS generation, highlighting the potential applications and future research directions in this field.

## 2.4 Talking Head

TH refers to a virtual or digital representation of the human face or head, typically used in multimedia applications, computer

TABLE 2: Multi-Modal Body Language Datasets.

| Type | Name | Year | Scale | Modal | Language | Link |
|---|---|---|---|---|---|---|
| Sign Language | Dicta-Sign [76] | 2008 | ∼1k | Video-Text | English | Link |
| | PHOENIX-Weather [77] | 2012 | ∼3k | Video-Text | Germany | Link |
| | ASLLVD [78] | 2012 | ∼3K | Video-Text | English | Link |
| | SIGNUM [79] | 2013 | ∼33K | Video-Text | Germany | Link |
| | DEVISIGN [80] | 2014 | ∼24k | Video-Text | Chinese | Link |
| | ASL-LEX 1.0 [81] | 2017 | ∼1K | Video-Text | English | Link |
| | PHOENIX14T [82] | 2018 | ∼68K | Video-Text | Germany | Link |
| | CMLR [83] | 2019 | ∼102K | Image-Text | Chinese | Link |
| | KETI [84] | 2019 | ∼15K | Video-Text | Korean | Not Available |
| | GSL [85] | 2020 | ∼3K | Video-Text | Greek | Link |
| | ASL-LEX 2.0 [86] | 2021 | ∼ 10K | Video-Text-Depth | English | Link |
| | How2sign [87] | 2021 | ∼35K | Video-Text-Skelton(2D)-Depth | English | Link |
| | Slovo [88] | 2023 | ∼20K | Video-Text | Russian | Link |
| | AASL [89] | 2023 | ∼8K | Image-Text | Arabic | Link |
| | ASL-27C [83] | 2023 | ∼23K | Image-Text | English | Link |
| Cued Speech | FCS [90] | 2018 | ∼13k | Video-Text-Audio | French | Link |
| | BEC [59] | 2019 | ∼3k | Video-Text-Audio | English | Link |
| | PCSC [91] | 2020 | 20 (P) | Video-Text-Audio | Polish | Link |
| | CLeLfPC [92] | 2022 | 350 | Video-Text-Audio | French | Link |
| | MCCS-2023 [41] | 2023 | ∼132k | Video-Text-Audio-Skelton(2,3D) | Chinese | Link |
| Co-speech | Trinity [67] | 2018 | 224(Min) | Videos-Text-Audio-Skelton(2,3D) | English | Link |
| | TED-Gesture [68] | 2019 | ∼252k | Videos-Text-Audio-Skelton(2D) | English | Link |
| | Talking With Hands [69] | 2019 | 200 | Videos-Text-Audio-Skelton(2,3D) | English | Link |
| | Speech2Gesture [70] | 2019 | ∼60k | Videos-Text-Audio-Skelton(2D) | English | Link |
| | TED-Expressive [27] | 2022 | ∼252k | Videos-Text-Audio-Skelton(2,3D) | English | Link |
| Talking Head | GRID [93] | 2006 | ∼34k | Video-Text | English | Link |
| | eNTERFACE [94] | 2006 | ∼1k | Video-Text-Audio | Multiple | Link |
| | MIRACL-VC1 [95] | 2014 | ∼3k | Video-Text-Depth | English | Link |
| | CREMA-D [96] | 2015 | ∼7k | Video-Text-Audio | English | Link |
| | TCD-TIMIT [96] | 2015 | ∼7k | Video-Text-Audio | English | Link |
| | MODALITY [97] | 2015 | ∼6k | Video-Text-Audio | English | Link |
| | LRW [98] | 2016 | ∼539k | Video-Text | English | Link |
| | MSP-IMPROV [99] | 2016 | ∼ 1K | Video-Text-Audio | English | Link |
| | ObamaSet [100] | 2017 | ∼1k | Video-Text-Audio | English | Link |
| | VoxCeleb1 [101] | 2017 | ∼22k | Video-Text-Audio | English | Link |
| | VoxCeleb2 [102] | 2018 | ∼146k | Video-Text-Audio | English | Link |
| | LRS2 [103] | 2018 | ∼96k | Video-Text | English | Link |
| | LRS3-TED [104] | 2018 | ∼119k | Video-Text | English | Link |
| | RAVDESS [105] | 2018 | ∼1k | Video-Text-Audio | English | Link |
| | MELD [106] | 2018 | ∼13k | Video-Text-Audio | English | Link |
| | AVSpeech [107] | 2018 | ∼150k | Video-Audio | Multiple | Link |
| | VOCASET [108] | 2019 | 480 | Video-Text-Audio-3DFace | English | Link |
| | LRW-1000 [109] | 2019 | ∼718K | Video-Text | Chinese | Link |
| | FaceForensics++ [110] | 2019 | ∼1k | Video-Text-Audio | English | Link |
| | MEAD [111] | 2020 | ∼281k | Video-Text-Audio | English | Link |
| | HDTF [112] | 2021 | ∼10k | Video-Text-Audio | English | Link |
| | AnimeCeleb [113] | 2022 | ∼2.4M | Video-Text-Audio-3DFace | English | Link |
| | VLRDT [114] | 2022 | ∼2k | Video-Text | Turkish | Link |
| | KoEBA [115] | 2023 | 104(P) | Video-Text-Audio | Korea | Link |
| Others | AV Letters [116] | 2002 | ∼19k | Video-Text | English | Link |
| | AV Digits [117] | 2002 | ∼5k | Video-Text | English | Link |
| | Aoyama Gakuin [118] | 2017 | ∼1k | Videos-Text-Audio-Skelton(2D) | Japanese | Not Available |
| | P2PSTORY [119] | 2018 | ∼13k | Video-Text-Audio | Multiple | Link |
| | AMASS [120] | 2019 | ∼18k | video-text-Skelton(3D) | English | Link |
| | BoLD [121] | 2020 | ∼10k | Video-Text-Audio-Skelton(3D) | English | Link |
| | PATS [122] | 2020 | ∼84k | Videos-Text-Audio-Skelton(2D) | English | Link |
| | BABEL [123] | 2021 | ∼28k | video-text-Skelton(3D) | English | Link |
| | HumanML3D [124] | 2022 | ∼15k | video-text-Skelton | English | Link |
| | BEAT [125] | 2023 | ∼3k | Video-Text-Audio-Skelton(3D) | Multiple | Link |

graphics, and HCI. It is usually an animated character that appears on a screen and can simulate various facial expressions, head actions, and speech with synchronized lip movements [17], [126], [127]. TH aims to enhance user experiences in various applications, from virtual assistants to entertainment platforms, by providing interactive and immersive communication interfaces.

In 2003, visual text-to-speech (VTTS) that generates talking faces driven by a speech synthesizer has been proposed for HCI systems [128]. Speech synthesis techniques are used to convert text input into synthesized speech, allowing the virtual character to speak. Facial animation algorithms are employed to animate the virtual character's facial movements, including lip synchronization with the generated speech. These algorithms analyze the phonetic information in the speech and map it onto corresponding facial movements. Additionally, sophisticated computer graphics techniques are utilized to generate realistic textures, lighting, and shading for the virtual character, enhancing its visual appearance [129]. Previous approaches to TH generation faced many limitations and were unable to achieve high-quality and realistic results due to constraints like limited data availability and computing power.

TH generation needs to fuse and synchronize information from different modalities to ensure consistency and coherence between the animation, sound, and text of the character. This involves the alignment, fusion, and synchronization of each modal to produce a more uniform response. In recent years, DL and multi-modal neural networks advance the performance of TH generation from multiple perspectives [19], [50]. By using multi-modal or cross-modal techniques based on a large amount of data, TH generation can integrate user input from different sources and interact in a more natural and realistic way. This enables multi-modal HCI systems to better understand user intentions, generate responses accordingly, and provide a more immersive and personalized interactive experience.

In Section 4.4, we give a brief review of TH recognition. In Section 5.4, we explore the applications and advancements in deep multi-modal learning for TH generation. We discuss the challenges associated with creating realistic and expressive virtual characters, including the synthesis of natural-sounding speech and the accurate representation of facial expressions. We review the SOTA techniques and highlight the potential future developments in this field, aiming to improve the realism and interactivity of THs in various applications.

# 3 BODY LANGUAGE DATASET AND EVALUATION METRICS

## 3.1 Body Language Datasets

Datasets have played a crucial role in the entire history of BL research, serving as a common foundation not only for measuring and comparing the performance of competing algorithms but also for driving the field toward increasingly complex and challenging problems. Particularly in recent years, DL techniques have brought significant success to BL research, with a substantial amount of annotated data being key to this success. The availability of large-scale image collections through the internet has made it possible to construct comprehensive datasets. Additionally, the availability of multi-modal data has provided richer information for related tasks, opening up new possibilities for future BL recognition and generation research.

In this section, we have collected and presented relevant datasets pertaining to BL tasks. As shown in Table 2, we have categorized them into five types based on data format and task purposes: CS, SL, CoS, TH, and others (Here "others" means these datasets are multi-modal BL datasets but are not designed for these four tasks). We have introduced the relevant information about these datasets, including publication year, dataset scale, available modalities, and the languages used in the datasets. Moreover, we have provided official links to these datasets to facilitate easier access for researchers. Please note that we measure the dataset scale based on the number of video clips/sequences. For datasets that do not provide these numbers, we provide the duration of the videos in minutes (represented as "Min") or the number of performers (represented as "P"). Some examples of BL datasets are shown in Figure 6 for more details.

We present the distribution of dataset languages in Figure 7. The chart shows the related datasets are primarily English datasets, but it also includes datasets in other languages like English datasets, Chinese datasets, and German datasets. This illustrates that current BL research is predominantly focused on English, but there is also growing importance placed on cross-cultural and multilingual datasets. Another problem is the difference in the format and standards of BL datasets. Different datasets may have varying storage formats, recording requirements, and model standards.

## 3.2 BL Generation Metrics

In order to evaluate the performance of BL gesture generation methods, we summarize the main gesture generation metrics and show these generation metrics, and corresponding calculation formulas in Table 6. A total of seven metrics are introduced for evaluation, namely PCK [130], FGD [131], MAE [132], STD [132], PMB [25], MAJE [131], and MAD [131]. Percentage of Correct Keypoints (PCK) assesses the accuracy of generated motion by comparing keypoints with actual motion. A predicted keypoint is considered correct if it falls within a specified threshold of the actual keypoints. Mean Absolute Error (MAE) quantifies the average difference between standardized coordinate values of generated and actual keypoints. Standard Deviation (STD) represents the variability or distribution of keypoints from their mean position after standardization. Fréchet Gesture Distance (FGD) measures dissimilarity between the distributions of latent features in generated and ground truth gestures, incorporating both location and spread. Percentage of Matched Beats (PMB) considers a motion beat matched if its temporal distance to an audio beat is below a threshold. Mean Absolute Joint Error (MAJE) calculates average errors between generated and ground truth joint positions across all time steps and joints. Mean Absolute Difference (MAD) computes average differences in joint accelerations, considering magnitude and direction. These criteria provide comprehensive insights into the accuracy, similarity, and alignment between generated and ground truth motion data.

The TH generation results can be evaluated quantitatively from multiple perspectives. Evaluation metrics include identity-preserving metrics, audio-visual synchronization metrics, image quality-preserving metrics, expression metrics, and eye-blinking metrics.
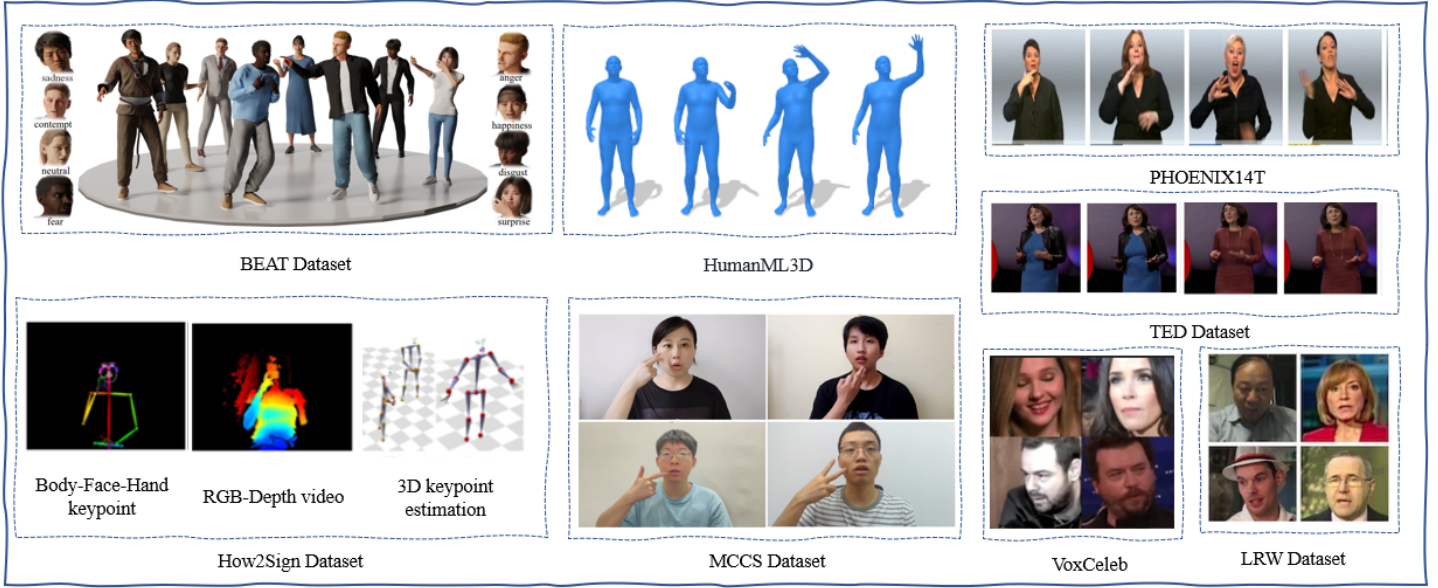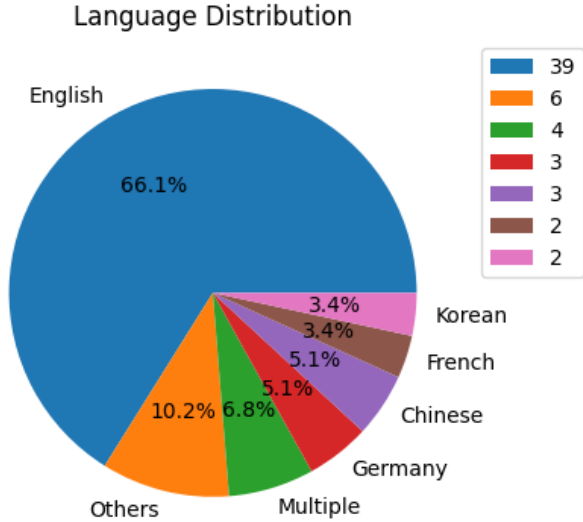
Fig. 6: Some examples of BL datasets.



Fig. 7: The Distribution of language used in BL datasets.

## 4 AUTOMATIC BODY LANGUAGE RECOGNITION

Here, we will introduce the recognition for the four BL variants, with a particular focus on the application expansion and innovation of multi-modal learning. In Figure 8, we present a summary of some representative works for BL recognition.

### 4.1 Sign Language Recognition

SLR aims to utilize a classifier to recognize SL glosses from video streams. It can be classified into two types in general according to the content of the SL: continuous SL recognition and isolated SL recognition. In this paper, we focus on continuous SL recognition (CSLR), in which the feature encoder module first extracts semantic representations from the sign video, and then the sequential module performs the mapping from the extracted semantics to the text sequence. In addition, some training strategies have been investigated for sufficient training. The comprehensive approaches are presented in Figure 9.

**Feature Encoder.** Since the hand acts a dominant role in the expression of SL, it has evolved over the past three decades, and we can divide these methods into the two following types.

- **Handcraft-based Method.** In the early research, handcrafted features are used to extract the hand motion, shape and orientation, such as HOG [77], [133], Grassmann covariance matrix (GCM) [134] and SIFT [135]. However, these methods require manual feature extraction and cannot directly be applied to different gestures, which means that different gestures necessitate distinct feature extraction approaches, resulting in a substantial amount of work.

- **CNN-based Method.** With the development of DL, CNNs [136], [137], [138], [139] generally replace the handcraft-based methods, becoming the most powerful feature extractor for SLR. Many researchers try to explore the reasonable CNN-based architecture to directly extract discriminative visual features from the video sequence. Specifically, exist works used 2D-CNN-TCN [44], [140], [141], [142] and 3D-CNN [140], [143], [144], [145], [146], [147], [148] as the backbone to extract spatial-temporal discriminative cues. For instance, IAN [143] utilizes 3D-ResNet [138] for visual representation. DNF [149] subtly designs 2D-CNN with the 1D temporal convolution, which has become one of the mainstream baseline methods. Although CNN-based methods can effectively capture spatial features in gesture images, they are limited in handling the temporal dynamics of gestures directly, and 3D-CNN-based methods involve significant computational overhead.

**Sequential Module.** There exist three representative approaches [150], [151] for CSLR. In the early research, HMM [72], [152], [153], [154] is used to learn the correspondence between the visual representation and sign gloss sequence. However, gesture actions of SLR often have long-term dependencies, and HMM
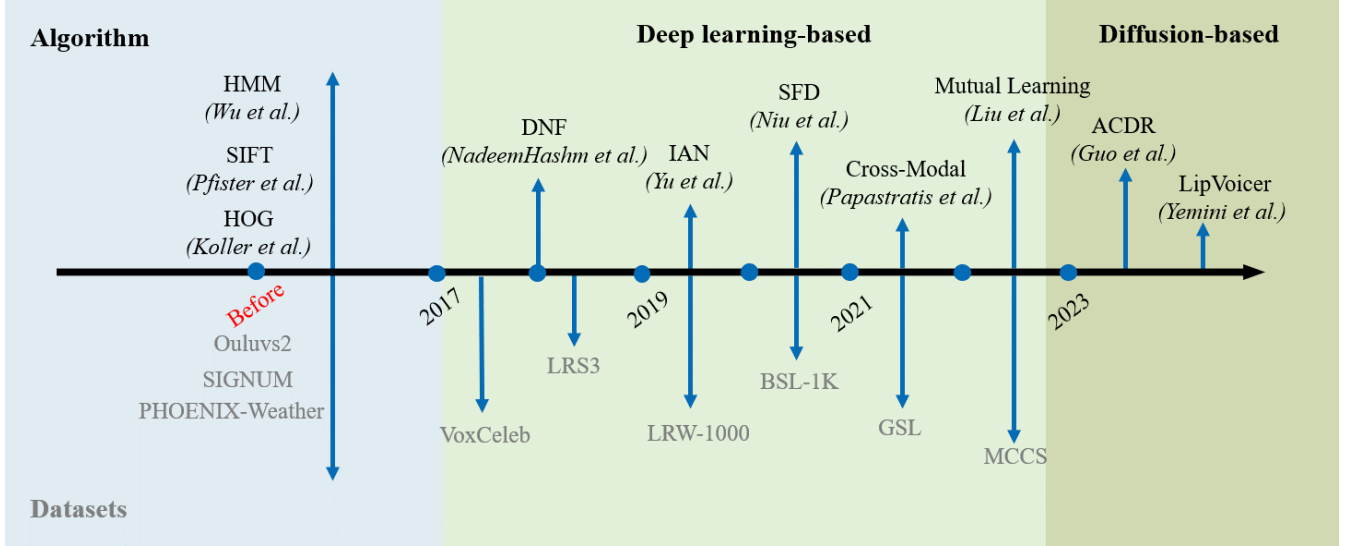
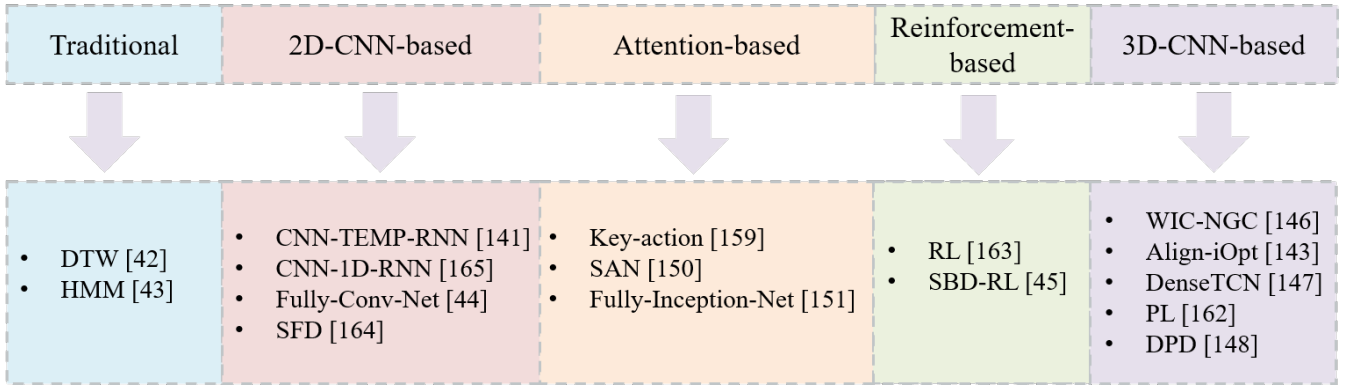Fig. 8: The milestones of Datasets and Methods for BL recognition.

Fig. 9: The comprehensive methods of SL recognition.

struggles to capture such complex sequential patterns. Additionally, HMM does not consider the alignment between input and output modalities. To this end, the RNN-based methods with CTC loss [141], [143], [155], [156], [157] are developed for CSLR to replace the HMM model, which improves the model's ability to handle data with incomplete alignments, but the ability to model global information is still limited. Therefore, to better understand the semantic relationship of the entire sign language sequence, encoder-encoder [22], [158], [159] has become a commonly used sequential framework. For instance, Guo et al. [22] utilizes the encoder-decoder framework with hierarchical deep recurrent fusion to merge cues from RGB and skeleton modalities.

**Training Strategy.** For sufficient training, some optimization strategies are widely used, with the most prominent being CTC [156], [160], [161] and Iterative Training [141], [142], [143], [161] strategies. On top of these two strategies, Pu et al. [142] introduce a cross-modality constraint called CMA to aid the training. Hao et al. [160] propose a three-stage optimization approach, which improves the recognition performance but it is time-consuming. Recently, Min et al. [156] further present two auxiliary constraints over the frame-level probability distributions, making the entire model end-to-end trainable.

## 4.2 Cued Speech Recognition

Automatic lip reading is a crucial component of ACSR. Therefore, we will first introduce the research progress in automatic lip-reading and then review ACSR.

### 4.2.1 Automatic Lip Reading

Advances in DL have led to a promising performance in lip-reading methods. Generally, DL-based lip-reading methods consist of two main parts, one is the extraction of visual feature information, and the other is the classification of sequence features.

**Feature Extraction.** Traditional studies use pixel-based [186], shape-based [187], [188], and hybrid-based [189], [190], [191], [192], [193] approaches to extract the visual feature. However, these methods are not only sensitive to image illumination change, lip deformation, and rotation but also cannot extract automatically.

Recently, DL has gradually become the mainstream research method in lip visual feature extraction, which can be divided into four categories. First, 2D-CNN-based methods are used [194], [195], which solves the problem of automatic feature extraction, but it can only process single-frame images and has a weak ability to process continuous frames, ignoring the spatio-temporal correlation between continuous frames. Then,

TABLE 3: The timeline of some representative works for BL recognition.

| Type | Year | Ref | Feature Extraction | Sequence Model | Learning Paradigm | Dataset |
|---|---|---|---|---|---|---|
| SL | 2019 | Pei et al. [162] | 3D-ResNet | BGRU | CTC | Phoenix-2014 |
| | 2019 | Pei et al. [163] | 3D-ResNet | Transformer | Reinforcement Learning | Phoenix-2014 |
| | 2019 | Cui et al. [141] | CNN | RNN | Iterative Training | Phoenix-2014 and SIGNUM |
| | 2020 | Niu et al. [164] | CNN | Transformer | CTC | Phoenix-2014 |
| | 2020 | SAFI [151] | 2D-CNN *plus* 1D-CNN | SAN | ACE *plus* CTC | Phoenix-weather and SIGNUM |
| | 2021 | Koishybay et al. [165] | 2D-CNN *plus* 1D-CNN | RNN | Iterative GR *plus* CTC | Phoenix-weather and SIGNUM |
| | 2021 | SLRGAN [166] | CNN | BiLSTM | GAN | Phoenix-weather, CSL and GSL |
| | 2022 | Chen et al. [167] | S3D | BLC | CTC *plus* Self-distillation | Phoenix-2014 |
| | 2022 | Zhou et al. [161] | SMC | BiLSTM *plus* SA-LSTM | CTC *plus* Keypoint Regression | PHOENIX-2014, CSL and PHOENIX-2014-T |
| | 2023 | Hu et al. [168] | 2D-CNN | 1D-CNN *plus* BiLSTM | SSTM *plus* TSEM | PHOENIX-2014, PHOENIX-2014-T, CSL and CSL-Daily |
| | 2023 | Zheng et al. [169] | CNN | VAE | CTC *plus* Contrastive Alignment Loss | PHOENIX-2014 and PHOENIX-2014-T |
| CS | 2018 | Liu et al. [90] | CNN | HMM | - | French CS |
| | 2021 | Papadimitriou et al. [170] | 2D-CNN *plus* 3D-CNN | Attention-based CNN | - | French and British English CS |
| | 2021 | Liu et al. [40] | CNN | MSHMM | HPM | French and British English CS |
| | 2021 | Wang et al. [58] | CNN *plus* ANN | BiLSTM *plus* FC | Cross-Modal Knowledge Distillation | French and British English CS |
| | 2022 | Sankar et al. [171] | Bi-GRU | Bi-GRU | CTC | CSF18 |
| | 2023 | Liu et al. [41] | ResNet-18 | Transformer | Cross-Modal Mutual Learning | Chinese, French, and British English CS |

3D-CNN-based methods have received extensive attention [196], [197], [198], [199]. Although this method can solve the problem of spatio-temporal correlation of continuous frames, it loses the extraction of fine-grained feature information by 2D convolution to a certain extent. According to the aforementioned issues, the hybrid methods [200], [201], [202] of 2D-CNN and 3D-CNN are also introduced to solve the problem of spatio-temporal feature extraction and local fine-grained feature extraction simultaneously. This method utilizes 3D-CNN to extract spatio-temporal information and then directly accesses 2D-CNN to extract fine-grained local information. However, it still affects the time information of feature coding to some extent. For that purpose, some other neural networks have gradually become a popular choice for lip visual feature extraction, such as Autoencoder model [203], [204], [205], [206].

**Recognition Modeling.** So far, there have been many works viewing lip reading as a sequence-to-sequence task and using sequence-based methods to deal with it, such as RNN, LSTM, and Transformer. It divides the feature representations extracted from the feature extractor into equal time steps, feeding each of them sequentially to the classification layer. For instance, [199], [202], [207], [208], [209], [210] utilize Long-Short Term Memory (LSTM) networks and Gated Recurrent Unit (GRU) to capture both global and local temporal information. Considering that Temporal Convolutional Network (TCN) has the advantage of faster converging speed with longer temporal memory than LSTM or RNN models, it is also widely used in this task. For example, Bai et al. [211] first propose a simple yet effective TCN architecture, indicating that TCN can become a reasonable alternative to RNN as a sequential model. Following this work, Martinez et al. [212] further demonstrate that multi-scale TCN

TABLE 4: The timeline of SL generation works.

| Type | Year | Ref | Input Modality | Framework | Dataset | Description |
|---|---|---|---|---|---|---|
| SL | 2011 | kippet et al. [172] | RCB video | EMBR | ViSiCAST | A gloss-centric tool is proposed to enable the comparison of avatars with human signers. But it is necessary to incorporate non-manual features. |
| | 2016 | John et al. [173] | RGB video | Segmental framework | Own dataset | This approach achieves automatic realism in generated images with low complexity, but it requires positioning the shoulder and torso. |
| | 2016 | Sign3D [174] | RGB video | Heterogeneous Database | Own dataset | This approach guarantees sign avatars that are easily understood and widely accepted by viewers, but it is restricted to a limited set of sign phrases. |
| | 2018 | HLSTM [22] | RCB video | LSTM | Own dataset | This approach shows robustness in effectively aligning the word order with visual content in sentences. Nevertheless, a limitation arises when generalizing it to new datasets. |
| | 2020 | Text2Sign [21] | Text | Transformer | PHOENIX14T | It demonstrates robustness in handling the dynamic length of the output sequence. However, It did not incorporate nonmanual information. |
| | 2020 | Zelinka et al. [175] | Text | CNN | Crech news | This method is robust to missing part, but face expression is not included. |
| | 2020 | ESN [23] | Text | GAN | PHOENIX14T | It shows Robustness to non-manual feature generation. But the genrated signs are not realistic . |
| | 2020 | Necati et al. [176] | Text | Transformers | PHOENIX14T | It does not need the gloss information, but the model is complex |
| | 2020 | Saunders et al. [177] | Text | GAN | PHOENIX14T | Robust to manual feature generation. The generated signs are not realistic. |
| | 2022 | DSM. [178] | Gloss | Transformer | PHOENIX14T | This work improves the prosody in generated Sign Languages by modeling intensification in a data-driven manner. |
| | 2022 | SignGAN. [179] | Text | FS-Net | meineDGS | It tackles large-scale SLP by learning to co-articulate between dictionary signs and improves the temporal alignment of interpolated dictionary signs to continuous signing sequences |
| | 2023 | PoseVQ-Diffusion. [180] | Gloss | CodeUnet | PHOENIX14T | It proposes a vector quantized diffusion method for conditional pose sequences generation and develops a novel sequential k-nearest-neighbors method to predict the variable lengths of pose sequences for corresponding gloss sequences |

TABLE 5: The timeline of Co-speech and Cued Speech generation works.

| Type | Year | Ref | Input Modality | Framework | Dataset | Description |
|---|---|---|---|---|---|---|
| CoS | 2015 | DCNF [181] | Text | FC network | DIAC | This work integrated speech text, prosody, and part-of-speech tags to generate co-verbal gestures using a combination of FC networks and a Conditional Random Field (CRF). |
| | 2019 | S2G [70] | RGB video | CNN | S2G | This work presents a method for generating gestures with audio speech, utilizing cross-modal translation and training on unlabeled videos. But it relies on noisy pseudo ground truth for training |
| | 2020 | StyleGestures [182] | RCB video | LSTM | Trinity | It achieves natural variations without manual annotation and allows control over gesture style while maintaining perceived naturalness. |
| | 2021 | A2G [183] | Text | CVAE | Trinity | This work employed a CVAE to generate diverse gestures from speech input and involved a one-to-many mapping of speech-to-gesture. |
| | 2021 | Text2Gestures [184] | Text | Transformer | MPI-EBEDB | Their approach employed Transformer-based encoders and decoders to generate sequential joint positions based on the text and previous pose. |
| | 2022 | ZeroEGGS [185] | Text | Variational Framework | Own dataset | A VAE-based framework is utilized to generate style-controllable CoS gestures and allowed for the generation of stylized gestures by conditioning on a zero-shot motion example |
| | 2022 | DiffGAN [24] | Text | Diffusion Model | PATS | An adversarial domain-adaptation approach is proposed to personalize the gestures of a speaker |
| | 2022 | RG [25] | Trinity and TED | QVAE | PHOENIX14T | This work introduces a novel CoS gesture synthesis method that effectively captures both rhythm and semantics. |
| CS | 1998 | Paul et al. [28] | Text | Template | Own dataset | Relying on manually selected keywords, low-context sentences, and pre-defined gesture templates. Its limitations included constrained expressiveness and increased manual effort. |
| | 2008 | Gérard et al. [29] | RGB video | Template | Own dataset | A post-processing algorithm was introduced to fine-tune synthesized hand gestures by addressing rotation, translation, and adaptation to new images. However, it relies on prior knowledge for adapting to new images. |

can outperform RNN in lip reading isolated words. However, these methods are relatively weak in modeling long-term dependencies and cannot directly capture long-term dependencies in sequences. Therefore, a new trend in the use of Transformer [213] for lip-reading tasks has emerged [37], [214].

Although the aforementioned methods achieve promising performance, they cannot solve the problem of inconsistency between the input and the output modality for lip reading. For that purpose, many advanced works are further developed in recent years, such as attention mechanisms [98], [197], [214], [215], [216], [217], [218] and contrastive learning [219].

### 4.2.2 Automatic Cued Speech Recognition

The literature on ACSR can be classified into three main categories: Multimodal Feature Extraction, Multimodal Fusion, and ACSR Modeling. We discuss them separately in this section and review the representative works of CS in Table 3

**Multi-modal Feature Extraction.** In the literature, there are several popular methods for CS feature extraction (*i.e.*, lips, hand position and hand shape).

- **Traditional Method.** It uses artificial markings to record lips and hands from video images [220], [221]. For example, Burger et al. [222] let the speaker wear black gloves to obtain accurate hand segmentation, while Noureddine et al. [40] placed blue marks on the speaker's fingers to obtain the coordinates of the fingers. However, both the speaker's clothing color and the background color can affect the accuracy of the hand segmentation.
- **CNN-based Method.** Recently, some CNN-based methods are utilized to get rid of artificial markings. For example, the CNN model is used in [40], [90], [170] to extract visual features from the regions of the lip and hand. On the basis of using the CNN model for the feature extraction of lips and hand shape, Liu et al. [40] further adopt the artificial neural network (ANN) to process the hand position feature. However, although CNN-based methods do not require artificial marks, their performances are limited by data scarcity.
- **Other DL-based Method.** Considering the data-hungry problem for multi-modal, some researchers try to introduce some advanced methods to solve this issue. For instance, Wang et al. [58] use lips, hand shape, and hand position to pre-train multi-modal feature extractor, using it for feature extraction of ACSR task. In addition, in another of their work [206], the three-stage multi-modal feature extraction model based on self-supervised contrastive learning and self-attention mechanism is proposed to model spatial and temporal features of CS hand shape, lips, and hand position.

**Multi-modal Fusion.** Most existing works in ACSR tend to direct concatenate the multi-modal feature flows, letting the model learn such features implicitly [90], [170], [206], [221]. For instance, [206], [221] utilize artificial marks to obtain regions of interest (ROIs) and directly concatenated features of lip and hand. MSHMM [90] merges different features by giving weights for different CS modalities. However, to the best of our knowledge, a critical issue in ACSR is the asynchrony between hand and lip articulations [40], [57], [223], while these researches mainly assume lip-hand movements are synchronous by default, ignoring the asynchronous issue.

Therefore, to tackle asynchronous modalities in the ACSR task, Liu et al. [40] propose to utilize the re-synchronization method to align the hand and lips features, which is realized by introducing the prior knowledge of the hand position and hand shape. Nevertheless, since the acquisition of prior knowledge depends on speakers and specific datasets, it is difficult to directly apply it to other languages. For that purpose, Liu et al. [41] further propose a Transformer-based cross-modal mutual learning framework for multi-modal feature fusion. The framework captures linguistic information by constructing a modality-invariant shared representation and uses this linguistic information to guide cross-modal information alignment. Recently, [224] proposes a novel Federated CS recognition (FedCSR) framework to train an model of CS recognition in the decentralized data scenario. Particularly, they design a mutual knowledge distillation fusion mechanism to maintain cross-modal semantic consistency of the CS multi-modalities, which learning a unified feature space for both speech and visual feature.

**ACSR Modeling.** ACSR aims to transcript visual cues of speech to text. In the early research, traditional statistical methods are used, which map sequences of hand-crafted features to phonemes using statistical models, such as HMM [220], [221] and HMM-GMM [40], [90]. However, such methods only consider the relationships between the current state and the previous one, which means that longer contextual information cannot be captured.

More recently, traditional DL-based methods (*i.e.*, CNN-based, LSTM-based) have been developed to alleviate the aforementioned problem. For instance, Sankar et al. [39] propose a novel RNN model trained with a Connectionist Temporal Classification (CTC) loss [225]. Papatimitriou et al. [170] propose a fully convolutional model with a time-depth separable block and attention-based decoder. However, such traditional DL-based methods still cannot capture long-time dependencies well, while it would be desirable to capture global dependency [213] over dynamic longer because of the context relationships of phonemes in long-time CS videos. For that purpose, Transformer-based methods [170] receive a lot of attention on the ACSR task in recent years. This kind of method achieves promising performance on the ACSR task, but it still requires powerful computing resources and a large dataset for training and parameter tuning.

Therefore, considering the existing corpus for ACSR is limited, some advanced methods such as cross-modal knowledge distillation method [58] and contrastive learning method [206], are also introduced to this task.

## 4.3 Co-speech Recognition

Although the existing research on CoS mainly focuses on the generation of CoS gestures, some scholars have shown that recognizing emotional expressions in CoS are crucial to this generation task. For example, Bock et al. [226] is the first to use the EmoGes corpus for emotion recognition in CoS gesture generation. Bhattacharya et al. [227] proposed to leverage the Mel-frequency cepstral coefficients and the text transcript computed from the input speech in separate encoders in our generator to learn the desired sentiments and the associated affective cues.

## 4.4 Talking Head Recognition

Since the development of the TH generation still has a long way to go, the focus of recent studies is not on TH recognition.

TH recognition is primarily treated as an evaluation metric for TH generation algorithms. However, humans have the ability to detect and identify a person from their face, even when there are changes in gender or facial expressions. However, it is difficult to build an automatic face recognition system. Therefore, the focus of TH recognition is primarily on capturing the essential facial attributes of the target speaker rather than full-fledged recognition of the speaker's identity. In the work proposed by Wen et al. [228], they classified the face identity to assess the performance of voice-based face reconstruction for known subjects. For unknown subjects, they used a gender classifier to evaluate the gender of the generated faces. Additionally, the feature distance, such as Cosine, $L_1$, and $L_2$ distances, between the target face and the generated face can be calculated to measure the accuracy of the generated face. To achieve this, a pre-trained face recognition model like FaceNet [229] or ArcFace [230] is employed as a feature extractor. The landmark distance (LMD) can also be measured as the disparity between the generated face and the real-world target face images.

TABLE 6: Metrics for gesture generation.

| Metrics | Calculation Formula |
|---------|---------------------|
| PCK [130] | $\text{PCK} = \frac{1}{N}\sum_{i=1}^{N}\mathbf{1}(d_i \leq \tau)N$ |
| FGD [131] | $\text{FGD} = \max_{\pi}\left(\frac{1}{T}\sum_{t=1}^{T}d(g_t^*, g_{\pi(t)})\right)$ |
| MAE [132] | $\text{MAE} = \frac{1}{T}\sum_{t=1}^{T}|g_t - g_t^*|$ |
| STD [132] | $\text{STD} = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(g_i - \bar{g})^2}$ |
| PMB [25] | $\text{PMB} = \frac{1}{N_m}\sum_{i=1}^{N_m}\sum_{j=1}^{N_a}\mathbf{1}\left[\left\|\boldsymbol{b}_i^m - \boldsymbol{b}_j^a\right\|_1 < \delta\right]$ |
| MAJE [131] | $\text{MAJE} = \frac{1}{N\cdot T}\sum_{t=1}^{T}\sum_{n=1}^{N}|g_t^n - g_t^{*n}|$ |
| MAD [131] | $\text{MAD} = \frac{1}{N\cdot T}\sum_{t=1}^{T}\sum_{n=1}^{N}|a_t^n - a_t^{*n}|_2$ |

*The corresponding meanings for the alphabet are as follow: N – The number of samples; $d_i$ – The distance of generated points and ground truth; $\tau$ – The threshold of PCK; T – the number of generated frames; $g_t$ – the generated gestures; $g_t$ – the ground truth; $b_i$ – The key frame corresponding to the beat. $a_i$ – The movement aceleration of generated gestures.*

## 5 AUTOMATIC BODY LANGUAGE GENERATION

The gesture generation task aims to generate a continuous sequence of gestures (*i.e.*, face, head, and hand) using multi-modal inputs (*e.g.*, gloss, speech, and text). In this section, we present the related works on gesture language generation and review the development timeline of gesture language generation applications, such as CS, SL, CoS gesture generations, and TH Generation, respectively.

### 5.1 Sign Language Generation

At the very beginning, we first present the difference between SL, CoS, and CS in Figure 4. SL generation has been studied for a long time. In this part, we mainly discuss the DL-based research on SL generation. For other SL generation methods, please refer to [5]. In Table 4, we present a summary of the details of the related SL generation works.

**Multi-modal Feature Extraction.** As a special visual language, the inputs of the SL gesture generation task are not only text and speech but also SL Gloss. It is a marking system for recording SL words and phrases, usually using written symbols and short descriptions to represent gestures, mouth movements as well as other non-gesture features. SL Gloss is suitable for recording the content of SL in written form to facilitate learners to learn and understand SL expressions. Previous work [177], [231] first converts spoken language to gloss and then uses gloss as input to extract features to generate SL gestures. Some work [175] use spoken language words and their characters as input to extract the word embedding of text, then the text features were used for gesture generation.

**Generative Methods.** For the SL generation task, there are several popular DL-based methods: 1) Neural Machine Translation (NMT) method, which [21], [22], [82] views the SL generation as a translation task. It uses the neural machine translation model to process SL text input, which can handle the output SL sequence of dynamic length but needs to solve problems such as domain adaptation. 2) Motion Graph method [21] uses motion graphics technology to construct a directed graph from motion capture data and generate SL. This method can handle the continuity of SL, but it requires large scales of data and another challenge is the scalability and computational complexity of the graph to select the best transitions. 3) Conditional generation methods such as Generative Adversarial Networks (GAN) and Variational Auto-Encoders (VAEs) are also employed to generate SL videos. A hybrid model, including a VAE and GAN combination, has been proposed for the generation of people performing SL [232], [233]. However, the problems such as model complexity and video quality need to be solved. 4) Other methods. In addition to the previous work, some research tries to introduce novel transformer-based model architectures for SLP. For example, [231] proposes a Progressive Transformers to generate continuous sign sequences from spoken language sentences. [234] combines a transformer with a Mixture Density Network (MDN) to manage the translation from text to skeletal pose. Although these works have brought performance improvements, the cost of the model complexity cannot be ignored.

Even though the SL generation has made some progress, some challenges in the CSL generation are still unsolved, *i.e.*, 1) The SL relies on facial expression to identify the specific meaning and avoid ambiguity. But few works consider facial expressions. 2) The scale of the SL gestures library is very large. According to the official Chinese SL dictionary, there are about 5600 kinds of frequently used SL gestures. Most of the dataset only covers a small portion of all gestures, for example, [235] builds a CSL dataset with 500 categories. The huge number of gestures brings a huge cost for the DL-based models to construct the mapping relationship.

### 5.2 Cued Speech Generation

As a lip-hand aided system, CS requires generating both lip and hand gestures simultaneously. Therefore, it is very important to extract multi-modal features such as speech features and text features. Among them, speech features have a strong correlation with lip movement. At the same time, text features play an important role in determining hand shape and position according to the coding system. As depicted in Figure 11, the generation of multi-modal CS hand gestures from audio-text
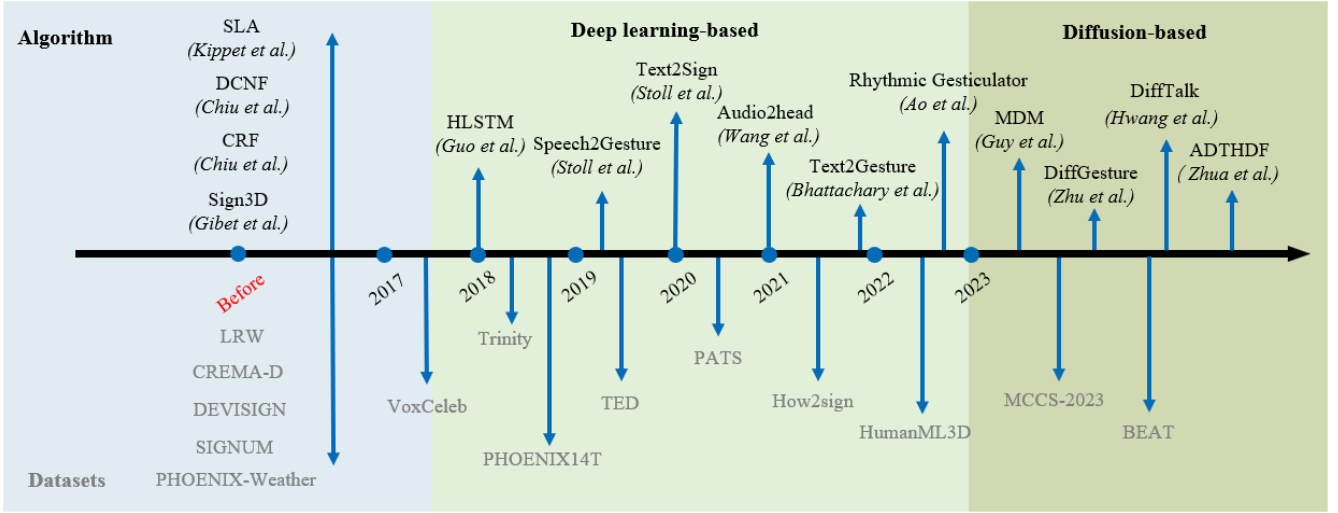
Fig. 10: The milestones of Datasets and Methods for BL generation.
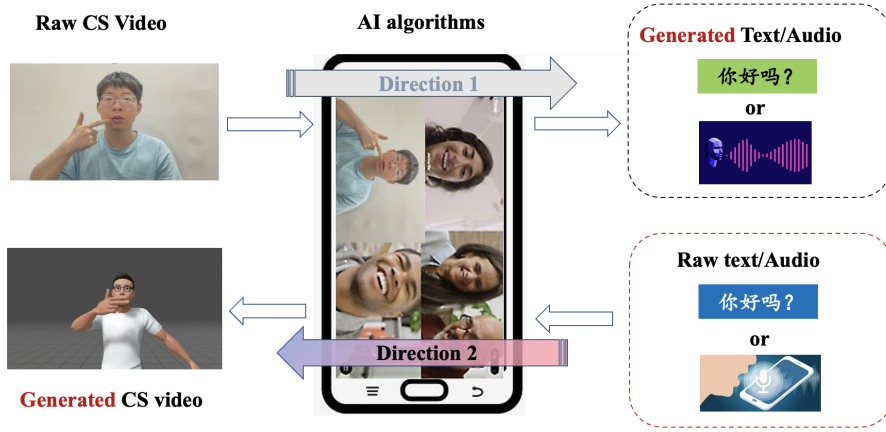


Fig. 11: The overall framework of the conversion between CS and text/audio. Direction 1 means CS to text/audio recognition, and direction 2 means text/audio to CS gesture generation. The first direction aims to recognize text or audio to make normal hearing better understand the hearing-impaired people, and the second direction can help the hearing-impaired to visually understand normal-hearing people.

is a crucial component of the CS conversion system. Previous studies in the literature have made limited initial attempts at CS gesture generation, which is mainly from two perspectives of multi-modal feature extraction and generation methods. Since the related work is relatively small, we incorporate the summary of the related CS generation works with the CoS In Table 5.

**Multi-modal Feature Extraction.** For CS generation, the feature includes continuous lip shape and hand shape movements. [28] used specific manually selected keywords, along with low-context sentences [236] as a feature, and pre-defined corresponding manual templates for hand gestures. CS recognition was performed, followed by the mapping of recognized text to the hand templates. However, this approach heavily relied on manual designs, which not only constrained the expressiveness of CS gestures but also increased the amount of manual effort required.

**Generative Method.** To the best of our knowledge, there is still a lack of research on end-to-end deep learning-based CS gesture generation. Only [28] proposed a post-processing algorithm to adjust synthesized hand gestures, involving correction of hand rotation and translation, as well as adaptation of the algorithm to new images. Nevertheless, this method requires prior human knowledge to adapt the algorithm to new images, leading to limited robustness.

## 5.3 Co-speech Generation

The milestones of CoS generation in recent years are presented in Figure 10. The upper part is related datasets and the lower part is the algorithm. The target of CoS gesture generation is to generate a sequence of body movements based on the corresponding audio input. It has been widely used in virtual character animation, especially in virtual speech and advertising. We divided it into three stages based on performance and popularity, Which are rule/statistical-based methods, DL-based methods, and Diffusion-based methods. In Table 5, we present a summary of the details of the related SL generation works.

**Multi-modal Feature Extraction.** In the CoS gesture generation task, the data of different modalities such as text and speech contain semantic and rhythmic information. How to extract and fuse these features to get a better representation is an important

topic. [131] uses a tri-modal encoder to encode text, speech, and person IDs separately, and then perform feature fusion, sampling from the fused feature space to complete the generation task. [237] separately models speech and text information. Instead of directly fusing at the feature level, it establishes two pipelines to model the dynamic and semantic information of the gesture motion, so as to generate accurate and rhythmic gesture sequences.

**Generative Model.** Numerous endeavors have been made in the process of choosing the generative model for CoS gesture generation task. In the early research, rule-based approaches [238], [239], [240] were used, which required the manual construction of a gesture library and the development of rules mapping from spoken language to gestures in the library. These methods had limited flexibility and required expert knowledge, but it is easier to be interpreted and were effective at handling semantic gestures. Then, statistical-based methods [241] replaced the manually written rules with traditional statistical models (*e.g.*, HMMs) trained on a dataset but still required the high-cost manual construction of a gesture library. In recent years, DL-based end-to-end approaches [25], [242] have been developed, which use raw "speech-gesture" datasets such as Trinity and TED [67], [68] to train deep neural networks for end-to-end gesture generation. These methods have reduced system complexity and produced more natural and fluid gestures, but they cannot guarantee the accuracy of generated rhythmic and semantic gestures. Meanwhile, most CoS research works do not consider the generation of the whole body, which also limits its expressiveness. Recently, diffusion models [243] have emerged as powerful deep generative models. Zhu et al. [244] introduced a novel diffusion-based framework called DiffGesture, which effectively captures the associations between audio and gestures and maintains temporal coherence to generate high-quality CoS gestures. However, the diffusion-based method has limitations in terms of training cost and the need for multiple steps to achieve satisfactory results, which hinders its real-time application in CoS gesture generation.

## 5.4 Talking Head Generation

TH generation has become an emerging research topic in recent years. As shown in Figure 12, talking face generation from an audio clip or dynamic TH generation from a target image and an audio clip are two fundamental research problems. The problems' solutions are essential to enabling a wide range of practical applications: (a) Entertainment: Generating virtual characters with realistic expressions and voice output can be applied to virtual reality games, special effects in movies, and other fields, to enhance user experience; (b) Virtual assistants: Generating virtual assistants with natural language voice and facial expressions can be used in customer service, robot assistants, and other scenarios to improve natural language interaction experience; (c) Human-machine interaction: Generating virtual characters with realistic expressions and voice output can be used for virtual meetings, remote education, and other scenarios to improve human-machine interaction effectiveness. (d) Healthcare: Generating virtual doctors with natural speak voices and facial expressions can be used in telemedicine, psychotherapy, and other scenarios to improve service quality and user experience.

### 5.4.1 Speech-to-face Generation

There is a strong connection between speech and face attributes, such as age, gender, and the shape of the mouth, which directly



(a) Overview of speech-to-face generation
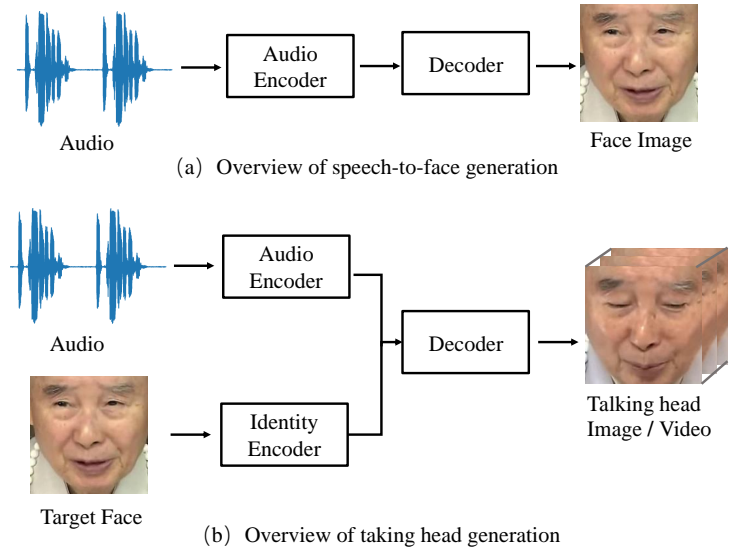


(b) Overview of taking head generation

Fig. 12: The two basic problems of speech-to-face generation.

affect the mechanics of speech generation [245]. Additionally, properties of speech such as language, accent, speed, and pronunciation are frequently shared among various nationalities and cultures. These properties can consequently manifest as standard physical facial features.

**Feature Extraction.** For speech and face feature extraction, Oh et al. [246] employs a trained face recognition network [247] to obtain face embedding, and a voice encoder that takes a complex spectrogram of speech as input and output speech features. Duarte et al. [248] design a speech encoder modified from SEGAN discriminator [249] to learn audio embedding. Similarly, Wen et al. [228] develops a voice embedding network consisting of six convolution layers to learn speech features. A voice encoder included voice activity detection and V-net is used in Fang et al. [250] to output audio embedding.

**Generation Model.** Oh et al. [246] employ a pre-trained face decoder [251] to reconstruct the face image. Motivated by the success of GAN [252] in generation images with high quality. Duarte et al. [248] developed a conditional GAN called WavPix that is able to generate face images directly from the speech. For better identity matching, Wen et al. [228] introduced the second discriminator to verify the identity of face image output. Considering that emotional expression is a key face attribute of a realistic face image, Fang et al. [250] applied two classifiers to measure identity and emotion semantic relevance in generating. In [253], a Face-based Residual Personalized Speech Synthesis Model (FR-PSS) containing a speech encoder, a speech synthesizer and a face encoder is designed for PSS.

These aforementioned methods can generate face images from speech, however, the authenticity and accuracy of the reconstructed face image still need to be improved: (a) Explicit cross-modal correlation learning is vital for identity information preservation, which is not explored in the previous methods. (b) The face images synthesized by the GAN-based or CNN-based generator lack details and authenticity.

**Evaluations Metrics.** For speech-to-face generation methods, identity information preservation is the key factor, therefore, quantitative metrics related to identity consistency are used to evaluate the performance, which includes landmark distance,

feature distance, and face attributes evaluation. Landmark distance is to calculate the distance of landmark (LMD) of generated face image and true face, where the landmark is achieved by Dlib [254] pre-trained DL methods such as FaceNet [229]. Feature distances are Cosine, $L_2$, and $L_1$ distances calculated between the feature of the true face and generated face. The face attributes are generally evaluated by attribution recognition accuracy like gender recognition, identity recognition, and face retrieval. The quality of generated face images is also important for speech-to-face generation, Fréchet Inception Distance (FID), and Inception score (IS) are two common metrics to evaluate performance. Those abovementioned metrics are highlighted in Table 7.

TABLE 7: Summary of quantitative metrics of Speech-to-face generation

| Metrics' degree | Metrics |
|---|---|
| Identity preservation | LDM [246], Cosine, $L_2$, $L_1$ [246], Face retrieval [246], Identity recognition [228], Gender classification [228] |
| Image quality | IS [250], FID [250] |

### 5.4.2 Talking Head Generation

Given a target face image and a speech clip, TH generation aims at synthesizing a sequence of target face images where the lip motion, head pose, and facial expressions are synchronized with the audio. Significantly different from the speech-to-face generation task, which extracts the identity of the speaker from the given speech, the TH generation task focuses on the content of the speech.

**Multi-modal Feature Extraction.** A VGG-M network pre-trained on the VGG Face dataset [293] is employed in [281] to learn face features and a speech encoder modified from VGG-M is used to learn speech embedding. Three temporal encoders are used to extract representations of the speaker's identity, the audio segment, and the facial expressions, and a polynomial fusion layer is designed to generate a joint representation of the three encodings [260]. Differently, Mittal et al. [266] develop a VAE to disentangle the phonetic content, emotional tone, and other factors into different representations solely from the input audio signal. To effectively disentangle each motion factor and achieve fine-grained controllable TH generation, Wang et al. [278] propose a progressive disentangled representation strategy by separating the factors in a coarse-to-fine manner, where we first extract unified motion feature from the driving signal, and then isolate each fine-grained motion from the unified feature. A pre-trained audio-to-AU module is employed in [270] to extract the speech-related AU information from speech.

**Multi-modal Learning.** For TH video generation, speech-synchronized lip movement, facial expressions, and head pose generation are key factors. Therefore, in the training stage, audio-visual cross-modal correlation learning is necessary for the consistency of these facial movements in a sequence. Chen et al. [255] propose an audio-visual correlation loss to synchronize lip changes and speech changes in a video regarding that variation along the temporal axis between two modalities are more likely correlated, specifically, the cosine similarity loss is used to maximize the correlation between the derivative of audio feature and

visual variations. For joint audio-visual representation learning, Zhou et al. [257] enforces the audio features and visual features to share a classifier so that they can share the same distribution, additionally, a contrastive loss is employed to close the paired audio and visual features. Eskimez et al. [263] designs a pair discriminator to improve the synchronization between the mouth shape and the input speech in the generated video. Zhu et al. [267] introduces the theory of mutual information neural estimation in talking face generation task to learn the cross-modal coherence.

**Generation Model.** The development of DL-based methods including CNN, RNN, GAN, Variational Autoencoder (VAE), Neural Radiance Fields (NeRF), and diffusion model (DM) have been explored in recent years. We compare the difference among them in Table 8 and Table 9.

The GAN-based methods are the mainstream for TH generation, in particular, because of their ability to synthesize data before the stronger generator DM emerged. In Table 8, we briefly list the recent works related to TH generation based on the GAN framework. Chen et al. [255] proposes a three-stream GAN to generate speech-synchronized lip video. Wang et al. [262] uses the GAN base network with an attentional mechanism to identify features related to head information. Zhang et al. [268] designs a FACIAL-GAN to encoder explicit and implicit attribute information for talking face video generation with audio-synchronized lip motion, personalized and natural head motion, and realistic eye blinks.

In addition to GAN-based approaches, inspired by the NeRF [294], Guo et al. [32] develops the audio-driven NeRF (AD-NeRF) model for TH synthesis, in which an implicit neural scene representation function is learned to map audio features to dynamic neural radiation fields for speaker face rendering. However, AD-NeRF often suffers from head and torso separation during the rendering stage. Therefore, a semantic-aware speaker portrait NeRF (SSP-NeRF) is proposed by Liu et al. [289]. They employ the semantic awareness of speech to address the problem of incongruity between local dynamics and global torso. The problem of slow rendering speed can not be ignored. To improve the real-time performance, Yao et al. [287] proposes a NeRF method that takes lip movement features and personalized attributes as two disentangled conditions, where lip movements are directly predicted from the audio inputs to achieve lip-synchronized generation.

Diffusion Probabilistic Models (DM) have shown strong ability in various generation tasks [295], [296]. Zhua et al. [291] proposes an audio-driven diffusion model for TH video generation, in which the lip motion features are aligned with the TH by contrastive learning. Yu et al. [290] proposes audio-to-visual diffusion prior trained on top of the mapping between audio and disentangled non-lip facial representations to semantically match the input audio while still maintaining both the photo-realism of audio-lip synchronization and the overall naturalness. Shen et al. [33] employs the emerging powerful diffusion models and model the TH generation as an audio-driven temporally coherent denoising process (DiffTalk). Xu et al. [292] first represents the emotion in the text prompt, which could inherit rich semantics from the CLIP, allowing flexible and generalized emotion control.

For better facial appearance transfer, intermediate faces such as 2D landmarks or 3DMM are widely used in TH generation. Figure. 13 illustrates a simplified pipeline of the TH generation methods based on intermediate face, which mainly consists of

TABLE 8: Summary of recent studies related to Talking Head generation. The following aspects are concluded: the network architecture for image synthesis and driving source; the methods work for a specific target or arbitrary identity; the audio feature is synchronized with lip motions or not; the ability to generate personalized attributes, and if any intermediate face models are used.

| Framework | Methods | Year | Driving source | Target | Audio features | Personalized | Face model |
|---|---|---|---|---|---|---|---|
| GAN | Chen et al. [255] | 2018 | Audio | Arbitrary | Sync | No | No |
| | Song et al. [256] | 2019 | Audio | Arbitrary | Sync | No | No |
| | Zhou et al. [257] | 2019 | Audio | Arbitrary | Sync | No | No |
| | ATVG [258] | 2019 | Audio | Arbitrary | not sync | No | 2D landmarks |
| | Vougioukas et al. [259] | 2019 | Audio | Arbitrary | Sync | Eye blinks, eyebrow | No |
| | Kefalas et al . [260] | 2020 | Audio | Arbitrary | No sync | No | No |
| | Sinha et al. [261] | 2020 | Audio | Arbitrary | No Sync | Eye blinking | No |
| | Wang et al. [262] | 2020 | Audio | Arbitrary | Sync | Head pose | 2D landmark |
| | Wav2lip [30] | 2020 | Audio | Arbitrary | Sync | No | No |
| | Eskimez et al. [263] | 2020 | Audio | Arbitrary | Sync | No | No |
| | Yi et al. [264] | 2020 | Video | Specific | Not sync | Head pose | 3DMM |
| | Chen et al. [265] | 2020 | Video | Arbitrary | Not sync | Head pose | 3DMM |
| | Mittal et al. [266] | 2021 | Audio | Arbitrary | Not sync | No | No |
| | MEAD [111] | 2020 | Audio | Arbitrary | Not sync | Emotion | No |
| | Zhu et al. [267] | 2021 | Audio | Arbitrary | No sync | No | No |
| | FACIAL [268] | 2021 | Video | Arbitrary | Not sync | Head pose, eye blinking | 3DMM |
| | Zhang et al. [112] | 2021 | Audio | Arbitrary | Sync | Head pose, eyebrow | 3DMM |
| | Si et al. [269] | 2021 | Audio | Arbitrary | No sync | Emotion | No |
| | Chen et al. [270] | 2021 | Audio | Arbitrary | Sync | No | No |
| | PC-AVS [271] | 2021 | Video | Arbitrary | Sync | Head pose | No |
| | GC-VAT [272] | 2022 | Video | Arbitrary | Sync | Head pose, expression | No |
| | Wang et al. [273] | 2022 | Audio | Arbitrary | Sync | Head pose | No |
| | EAMM [274] | 2022 | Video | Arbitrary | No sync | Emotion | No |
| | SPACE [275] | 2022 | Audio | Arbitrary | No sync | Head pose, emotion | 2D landmark |
| | DIRFA [276] | 2023 | Audio | Arbitrary | Sync | No | No |
| | DisCoHead [115] | 2023 | Video | Arbitrary | Sync | Head pose, eye blinking, eyebrow | No |
| | OPT [277] | 2023 | Audio | Arbitrary | No sync | Head pose, expression | 3DMM |
| | Wang et al. [278] | 2023 | Audio | Abitrary | Sync | Head pose, expression,gaze, eye blinking | No |
| | Zhang et al. [279] | 2023 | Audio | Abitrary | No sync | No | No |

two steps: low-dimensional driving source data are mapped into facial parameters; then rendering network is used to convert the learned facial parameters into high-dimensional video output.

**Evaluation Metrics.** Various perspectives reveal that the generated text-to-speech (TTS) output lacks the authenticity of human speech: (a) The target individual's face should match that of the synthetic video's speaker, (b) The generated speaker's mouth should synchronize the audio, (c) The produced TH video should be of a good caliber, (d) The expression of the speaker in the generated video should be natural and match the emotion of the audio, and (e) Eye blinking should be expected when talking.

Thus, the quantitative metrics of TH generation can be classified from these five views, as shown in Table 10.

**Audio Input Pre-processing.** Most of the TH generation works are audio signal driven. Here, we will introduce how previous work has dealt with speech signals in this field. In general, the audio waveform is resampled at 16KHz, and then the audio feature is computed [305]. Spectrogram, MFCC, and Fbank are the three mostly used audio features. Fang et al. [250] performs an ablation experiment on these three audio features, and they found that Fbank achieved the best performance, while the Spectrogram performed the worst FID. The reasons they

TABLE 9: Summary of recent studies related to Talking Head generation. The following aspects are concluded: The network architecture for image synthesis; Driving source; The methods work for a specific target or arbitrary identity; The audio feature is synchronized with lip motions or not; The ability to generate personalized attributes, and if any intermediate face models are used.

| Framework | Methods | Year | Driving source | Target | Audio features | Personalized | Face model |
|---|---|---|---|---|---|---|---|
| CNN | X2Face [280] | 2018 | Audio, video | Arbitrary | Sync | Head pose, expression | No |
| | Jamaludin et al. [281] | 2019 | Audio | Arbitrary | Sync | No | No |
| | Wen et al. [282] | 2020 | Video, audio | Arbitrary | No sync | Head pose, expression | 3DMM |
| | LipSync3D [283] | 2021 | Video | Specific | No sync | No | 3DMM |
| | Audio2head [31] | 2021 | Audio | Arbitrary | Sync | Head pose | 2D landmark |
| | Lu et al. [284] | 2021 | Audio | Specific | No sync | Head pose, eyebrow | |
| RNN | Bigioi et al. [285] | 2022 | Video, audio | Arbitrary | No sync | Head pose | 2D landmark. |
| VAE | SadTalker [286] | 2023 | Audio | Abitrary | No sync | Head pose, eye blinking | 3DMM |
| NeRF | AD-NeRF [32] | 2021 | Audio | Specific | No sync | No | No |
| | DFA-NERF [287] | 2022 | Video | specific | Sync | Eye blinking, head pose | No |
| | DFRF [288] | 2022 | Audio | Arbitrary | No sync | No | 3DMM |
| | SSP-NeRF [289] | 2022 | Video | Arbitrary | No sync | No | 3DMM |
| DM | Yu et al. [290] | 2022 | Audio | Arbitrary | Sync | Facial motion | No |
| | Zhua et al. [291] | 2023 | Video | Arbitrary | Sync | Eye blinking, head pose | 3DMM |
| | DiffTalk [33] | 2023 | Audio | Arbitrary | Sync | No | No |
| | Xu et al. [292] | 2023 | Audio, text | Arbitrary | No sync | Emotion | 3DMM |



Fig. 13: The typical pipeline of TH generation methods based on the intermediate face.

TABLE 10: Summary of quantitative metrics of Talking Head Generation

| Metrics' degree | Metrics |
|---|---|
| Identity-preserving | PSNR, SSIM [297], FID, LMD, LPIPS [298], CSIM, IS, ACD [299] |
| Audio-visual synchronization | AV Conf, AV Off [300], WER [259], $LMD_m$ [272], [299], $Sync_{conf}$ [272], LRSD, LRA [98] |
| Image quality preserving | CPBD [301], FDBM [302] |
| Expression | Classification accuracy [303] |
| Eye blinking | EAR [304], Blink rate, Blink median duration [268] |

guessed that Spectrogram contained much redundant information, MFCC discarded some related information, and Fbank kept balance. However, MFCC is used the most in the talking face generation.

# 6 CHALLENGES OF BL RECOGNITION AND GENERATION

The existing BL recognition and generation methods have not been capable of meeting real-world requirements under exposure to various challenges. In order to fully demonstrate the typical challenges of BL recognition and generation in the field

of BL, we elaborate in detail on SL, CS, and TH from three aspects: **subtasks challenges**, **datasets challenges** and **evaluation metrics challenges**.

## 6.1 Subtasks Challenges

To more fully illustrate the challenges of BL recognition and generation tasks, we split each major task into three subtasks, *i.e.*, Lip reading, SL recognition, and CS recognition. From the perspective of task definition, the TH task itself is more focused on the generation process. Moreover, limited by the development of the existing TH generation, it is difficult for researchers to capture the basic facial attributes of the target speaker. The

existing studies lack an exploration of TH recognition, so the challenge of TH recognition is not included in the discussion of subtask challenges in this survey. Current research on CoS predominantly concentrates on CoS gesture generation. While some studies have demonstrated a positive impact on the CoS Generation task, the majority of recent works do not prioritize CoS Recognition as a primary focus. So the challenge of CoS Recognition is not included in the discussion of subtask challenges in this survey.

The challenges of BL recognition tasks are mainly due to the efficiency of the cross-modal feature fusion, and the specific challenges of each subtask are as follows.

- **Lip Reading.** There are two primary challenges in automatic lip reading: intra-class difference and inter-class similarity. The former is hindered by factors such as speech emotion, speed, gender, age, skin color, and speech habits, making it difficult to distinguish variations within the same word category. Additionally, the semantic disparities between words used in different contexts significantly impact lip reading. The latter challenge stems from the abundance of word categories, leading to challenges in visually distinguishing similar-looking words belonging to different classes. Addressing these challenges is crucial for improving the accuracy and effectiveness of lip reading recognition systems, which have valuable implications for aiding communication for individuals with hearing impairments and advancing the field's applications.
- **Sign Language Recognition.** SL recognition encounters significant challenges arising from the pronounced variations in gestures, which seriously impede its accuracy. Moreover, factors like hand shape, illumination conditions, and resolution play pivotal roles in limiting SL recognition performance. Additionally, occlusion, including self-occlusion between fingers and occlusion between hands and other body parts, adversely affects feature fusion, becoming a key influencing factor in SL recognition. Another pressing challenge is the development of a real-time multilingual SL recognition system. Addressing these complexities is essential to advance the field and improve the efficiency and inclusivity of SL recognition technologies.
- **Cued Speech Recognition.** The primary obstacle in CS recognition is the hand preceding phenomenon [41], where the hand movements often occur faster than the corresponding lip movements, anticipating the next phoneme. This phenomenon hampers the efficiency of lip and hand feature fusion in CS recognition. Besides, due to variations in individual CS coding habits and styles, adaptability in multi-cuer scenarios is also a challenge.

The challenge of BL generation mainly stems from the stability and quality of the generated gesture, and the specific challenges of each subtask are as follows.

- **Cued Speech Generation.** In conclusion, the CS generation faces several challenges that need to be addressed for the development of effective systems. The lack of large-scale annotated datasets, the complexity of modeling CS gestures, and the need for accurate asynchronous alignment between cued signs and spoken words are

key challenges. Additionally, integrating audio and visual modalities and achieving generalization to new speakers and languages are important considerations. Overcoming these challenges through advancements in modeling ability, multi-modal fusion, and the availability of diverse datasets will contribute to the improvement of CS generation systems.
- **Sign Language Generation.** In the realm of SL production, numerous obstacles warrant attention, chief among them being domain adaptation and model collapse. The former obstacle arises from the inherent variations in word styles and meanings across different languages, necessitating effective adaptation strategies. Furthermore, a noteworthy challenge lies in the limited proficiency of generating uncommon and unseen words, hindering the overall performance of the system. Moreover, the persisting issues of model collapse, non-convergence, and instability within generative models further compound the complexities faced in Sign Language production. Addressing these multifaceted challenges is crucial for advancing the SOTA in this domain and facilitating more reliable and robust SL generation.
- **Co-speech Generation.** The generation process of CoS encounters challenges due to the presence of highly idiosyncratic and non-periodic spontaneous gestures. The accurate capture of finger motion poses difficulties, resulting in the manifestation of idiosyncratic gestures. Furthermore, the non-periodic nature of gestures arises from the substantial variation in gesture behavior.
- **Talking Head Generation.** TH generation confronts two primary challenges: information coupling and diversity targets. The former encompasses the synchronization of multiple facial elements, such as head posture, facial expression, lip movement, and background motion, while also addressing the "uncanny valley effect" [306], a phenomenon common in face generation where generated faces appear almost human-like but lack true realism, leading to discomfort. The latter challenge pertains to harmonizing temporal resolution and speech features across diverse data modalities, along with the complexity of defining visual quality as a clear training objective. Overcoming these challenges is crucial for advancing the field and achieving a more realistic and visually coherent TH generation.

## 6.2 Datasets Challenges

The current datasets for SL recognition and generation encounter significant limitations due to the high costs associated with data collection and manual annotation. This results in datasets with small-scale and weak annotations, hindering the progress of BL-related tasks. To create BL datasets, collaboration between language experts and native speakers is essential, further adding to the complexities and expenses involved. A potential solution to address these challenges is to explore self-supervised learning using unlabeled BL data [307], which could alleviate the need for extensive manual annotation.

Moreover, privacy protection poses another hurdle, as some large BL datasets [308], [309] are not publicly accessible. In light of the high costs and privacy concerns, a viable approach is to leverage existing wild online videos to collect the necessary

BL data. Similar to the training datasets used for Contrastive Language-Image Pre-training (CLIP) [310] and DALL-E [311], employing very large datasets can enhance the generalization capabilities of BL recognition and generation models.

Apart from the dataset challenges, the high costs associated with collecting and annotating 3D data contribute to the scarcity of large-scale 3D BL datasets. Consequently, the development of 3D BL Generation faces significant obstacles in understanding and processing 3D BL data effectively. Overcoming these challenges is essential to advance the field of BL recognition and generation, allowing for more efficient and accurate communication support for individuals with hearing impairments.

## 6.3 Evaluation Metrics Challenges

The primary nature of the BL recognition task lies in its classification essence, where simple and efficient classification accuracy serves as the prevalent evaluation metric. However, this paper shifts its focus to the BL generation task and the challenges it poses in terms of evaluation metrics. Subjective metrics utilized in the BL generation task prove to be costly, time-consuming and lack scalability. Metrics like human likeness and gesture appropriateness, although valuable, suffer from non-replicability and instability issues. On the other hand, objective metrics such as PSNR, SSIM, FID [312] and LRSD [17] offer advantages over subjective ones but come with limitations in assessing the similarity between gesture and speech, as well as the semantic appropriateness of gestures. Notably, unlike subjective metrics that evaluate human likeness, the existing literature rarely quantifies objective metrics measuring gesture diversity or various motion appropriateness aspects. These challenges highlight the need for robust and comprehensive evaluation metrics in the BL generation domain to ensure an accurate and meaningful assessment of generated Sign Language outputs.

## 7 FUTURE DISCUSSIONS

Through an extensive summary and analysis of the existing literature, this survey offers the following **discussions** and **new insights**:

1) The integration of large-scale multi-modal BL datasets and the establishment of a unified low-loss data format are key factors in advancing BL recognition and generation tasks. By collecting extensive datasets from diverse online videos, we can enhance the generalization and robustness of BL recognition and generation models for real-world scenarios. Additionally, the adoption of a unified data standard and adaptable conversion method allows for the seamless integration of different datasets and facilitates collaboration among researchers. This promotes interoperability between models, enabling efficient sharing and utilization of resources within the research community.

2) Recently, large-scale pre-training models such as Chat-GPT have achieved outstanding performance in various visual-linguistic cross-modal tasks. For instance, CLIP and various variations of the multi-modal CLIP model have emerged. However, they have the following drawbacks: a) they might not deeply connect different types of data as effectively as specialized models; b) they demand in terms of computing power due to their size. c) This

model might not allow fine-tuning for specific tasks and could struggle with specialized knowledge; d) it needs a lot of diverse data to work well and could be hard to interpret. To this end, how to build a large-scale multi-modal model for BL recognition and generation is a promising topic.

3) Besides, it was found that the ability of existing large-scale pre-training models to learn fine-grain features still needs to be improved [313]. In BL, fine-grained feature learning is essential, For example, hand positions and lip movements in CS and CoS needed to be accurately recognized and generated to ensure clarity and avoid ambiguity. Therefore, fine-grained BL recognition and generation is a feasible direction to improve their performance.

4) The multi-modal models in the task of BL recognition and generation are very susceptible to the perturbations (attacks) of different modalities, resulting in serious performance degradation. How to pre-train a robust and secure multimodal large-scale model for BL recognition and generation is an urgent problem to be solved.

5) An essential requirement for BL recognition and generation systems is real-time capability, especially for multilingual and multiple-speakers scenarios. Creating a real-time system is vital to cater to the needs of both the deaf and speaking communities. However, existing audio-visual datasets are predominantly monolingual, with English being the most commonly represented language. In practical applications, multilingual communication is often necessary, highlighting the need for diverse datasets. Additionally, current methods for BL recognition and generation are often limited to specific target identities, as different speakers exhibit significant variations in appearance and habits. Overcoming these challenges is crucial to develop adaptable and effective real-time BL systems that accommodate various languages and diverse speakers.

## 8 CONCLUSION

This survey has delved into the realm of deep multi-modal learning for automatic BL recognition and generation, shedding light on its potential and challenges. This survey focuses on four classical BL variants, *i.e.*, Sign Language, Cued Speech, Co-speech, and Talking Head. Through a meticulous examination of various modalities, including visual, auditory, and textual data, and their integration, we have explored the intricacies of capturing and interpreting these four BL. By reviewing SOTA methodologies, such as feature fusion, representation learning, recognition, and generation methods, we have uncovered the strengths and limitations of current approaches. The significance of datasets and benchmarks in facilitating research progress was also emphasized, with a focus on annotation methodologies and evaluation metrics. Despite the progress, challenges persist, demanding the creation of diverse datasets, addressing limited labeled data, enhancing model interpretability, and ensuring robustness across environments and cultural contexts. Looking ahead, the future holds promises of more sophisticated architectures and training strategies, harnessing the complementary nature of multi-modal data and leveraging advancements in multi-modal learning,

large-scale pre-trained model, self-supervised learning, and reinforcement learning. As this research area evolves, it is poised to revolutionize human-human and human-machine interactions, fostering natural and effective communication across domains.

# REFERENCES

[1] R. O. Cornett, "Cued speech," *American annals of the deaf*, vol. 112, no. 1, pp. 3–13, 1967.

[2] B. Joksimoski, E. Zdravevski, P. Lameski, I. M. Pires, F. J. Melero, T. P. Martinez, N. M. Garcia, M. Mihajlov, I. Chorbev, and V. Trajkovik, "Technological solutions for sign language recognition: A scoping review of research trends, challenges, and opportunities," *IEEE Access*, vol. 10, pp. 40 979–40 998, 2022.

[3] X. Liu, Q. Wu, H. Zhou, Y. Du, W. Wu, D. Lin, and Z. Liu, "Audio-driven co-speech gesture video generation," *Advances in Neural Information Processing Systems (NIPS)*, vol. 35, pp. 21 386–21 399, 2022.

[4] B. Zhang, C. Qi, P. Zhang, B. Zhang, H. Wu, D. Chen, Q. Chen, Y. Wang, and F. Wen, "Metaportrait: Identity-preserving talking head generation with fast personalized adaptation," in *Proc. IEEE/CVF-CVRP*, 2023, p. 22096–22105.

[5] R. Rastgoo, K. Kiani, S. Escalera, and M. Sabokrou, "Sign language production: A review," in *Proc. IEEE/CVF-CVRP*, 2021, pp. 3451–3461.

[6] S. Nyatsanga, T. Kucherenko, C. Ahuja, G. E. Henter, and M. Neff, "A comprehensive review of data-driven co-speech gesture generation," in *Computer Graphics Forum*, vol. 42, no. 2. Wiley Online Library, 2023, pp. 569–596.

[7] R. Rastgoo, K. Kiani, and S. Escalera, "Sign language recognition: A deep survey," *Expert Systems with Applications*, vol. 164, p. 113794, 2021.

[8] A. Fernandez-Lopez and F. M. Sukno, "Survey on automatic lip-reading in the era of deep learning," *Image and Vision Computing*, vol. 78, pp. 53–72, 2018.

[9] S. Fenghour, D. Chen, K. Guo, B. Li, and P. Xiao, "Deep learning-based automated lip-reading: A survey," *IEEE Access*, vol. 9, pp. 121 184–121 205, 2021.

[10] R. Chand, P. Jain, A. Mathur, S. Raj, and P. Kanikar, "Survey on visual speech recognition using deep learning techniques," in *Proc. IEEE-CSCITA*, 2023, pp. 72–77.

[11] S. Bhaskar, T. Thasleema, and R. Rajesh, "A survey on different visual speech recognition techniques," in *Data Analytics and Learning (DAL)*, 2018, pp. 307–316.

[12] N. Radha, A. Shahina *et al.*, "A survey on visual speech recognition approaches," in *Proc. IEEE-ICAIS*, 2021, pp. 934–939.

[13] O. Koller, "Quantitative survey of the state of the art in sign language recognition," *arXiv preprint arXiv:2008.09918*, 2020.

[14] I. Adeyanju, O. Bello, and M. Adegboye, "Machine learning methods for sign language recognition: A critical review and analysis," *Intelligent Systems with Applications*, vol. 12, p. 200056, 2021.

[15] I. Papastratis, C. Chatzikonstantinou, D. Konstantinidis, K. Dimitropoulos, and P. Daras, "Artificial intelligence technologies for sign language," *Sensors*, vol. 21, no. 17, p. 5843, 2021.

[16] D. M. Madhiarasan, P. Roy, and P. Pratim, "A comprehensive review of sign language recognition: Different types, modalities, and datasets," *arXiv preprint arXiv:2204.03328*, 2022.

[17] 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.

[18] T. Sha, W. Zhang, T. Shen, Z. Li, and T. Mei, "Deep person generation: A survey from the perspective of face, pose, and cloth synthesis," *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–37, 2023.

[19] R. Zhen, W. Song, Q. He, J. Cao, L. Shi, and J. Luo, "Human-computer interaction system: A survey of talking-head generation," *Electronics*, vol. 12, no. 1, p. 218, 2023.

[20] C. Sheng, G. Kuang, L. Bai, C. Hou, Y. Guo, X. Xu, M. Pietikäinen, and L. Liu, "Deep learning for visual speech analysis: A survey," *arXiv preprint arXiv:2205.10839*, 2022.

[21] S. Stoll, N. C. Camgoz, S. Hadfield, and R. Bowden, "Text2sign: Towards sign language production using neural machine translation and generative adversarial networks," *International Journal of Computer Vision*, vol. 128, pp. 891–908, 2020.

[22] D. Guo, W. Zhou, H. Li, and M. Wang, "Hierarchical lstm for sign language translation," in *Proc. Conf AAAI Artif. Intell.*, vol. 32, no. 1, 2018.

[23] B. Saunders, N. C. Camgoz, and R. Bowden, "Everybody sign now: Translating spoken language to photo realistic sign language video," *arXiv preprint arXiv:2011.09846*, 2020.

[24] C. Ahuja, D. W. Lee, and L.-P. Morency, "Low-resource adaptation for personalized co-speech gesture generation," in *Proc. IEEE/CVF-CVPR*, June 2022, pp. 20 566–20 576.

[25] T. Ao, Q. Gao, Y. Lou, B. Chen, and L. Liu, "Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings," *ACM Transactions on Graphics (TOG)*, vol. 41, no. 6, p. 1–19, 2022.

[26] Y. Liang, Q. Feng, L. Zhu, L. Hu, P. Pan, and Y. Yang, "Seeg: Semantic energized co-speech gesture generation," in *Proc. IEEE/CVF-CVPR*, June 2022, pp. 10 473–10 482.

[27] X. Liu, Q. Wu, H. Zhou, Y. Xu, R. Qian, X. Lin, X. Zhou, W. Wu, B. Dai, and B. Zhou, "Learning hierarchical cross-modal association for co-speech gesture generation," in *Proc. IEEE/CVF-CVPR*, 2022, pp. 10 462–10 472.

[28] P. Duchnowski, L. D. Braida, D. Lum, M. Sexton, J. Krause, and S. Banthia, "Automatic generation of cued speech for the deaf: status and outlook," in *International Conference on Auditory-Visual Speech Processing (AVSP)*, 1998.

[29] G. Bailly, Y. Fang, F. Elisei, and D. Beautemps, "Retargeting cued speech hand gestures for different talking heads and speakers," in *Retargeting cued speech hand gestures for different talking heads and speakers*, September 2008, p. 8.

[30] P. KR, M. Rudrabha, P. Namboodir, and C. Jawahar, "A lip sync expert is all you need for speech to lip generation in the wild," in *Proc. ACM MM*, 2020.

[31] S. Wang, L. Li, Y. Ding, C. Fan, and X. Yu, "Audio2head: Audio-driven one-shot talking-head generation with natural head motion," in *Proc. IJCAI*, 2021.

[32] Y. Guo, K. Chen, S. Liang, Y.-J. Liu, H. Bao, and J. Zhang, "Ad-nerf: Audio driven neural radiance fields for talking head synthesis," in *Proc. IEEE/CVF-ICCV*, 2021, pp. 5784–5794.

[33] S. Shen, W. Zhao, Z. Meng, W. Li, Z. Zhu, J. Zhou, and J. Lu, "Difftalk: Crafting diffusion models for generalized audio-driven portraits animation," in *Proc. IEEE/CVF-CVPR*, 2023, pp. 1982–1991.

[34] P. Lucey, G. Potamianos, and S. Sridharan, "Patch-based analysis of visual speech from multiple views," in *International Conference on Auditory-Visual Speech Processing (AVSP)*. AVISA, 2008, pp. 69–74.

[35] Z. Zhou, G. Zhao, and M. Pietikäinen, "Towards a practical lipreading system," in *Proc. IEEE/CVF-CVPR*, 2011, pp. 137–144.

[36] P. Wu, H. Liu, X. Li, T. Fan, and X. Zhang, "A novel lip descriptor for audio-visual keyword spotting based on adaptive decision fusion," *IEEE Transactions on Multimedia*, vol. 18, no. 3, pp. 326–338, 2016.

[37] P. Ma, S. Petridis, and M. Pantic, "End-to-end audio-visual speech recognition with conformers," in *Proc. IEEE-ICASSP*, 2021, pp. 7613–7617.

[38] L. Liu, G. Feng, D. Beautemps, and X.-P. Zhang, "A novel resynchronization procedure for hand-lips fusion applied to continuous french cued speech recognition," in *Proc. IEEE-EUSIPCO*, 2019, pp. 1–5.

[39] K. Papadimitriou, M. Parelli, G. Sapountzaki, G. Pavlakos, P. Maragos, and G. Potamianos, "Multimodal fusion and sequence learning for cued speech recognition from videos," in *International Conference on Human-Computer Interaction*, 2021, pp. 277–290.

[40] L. Liu, G. Feng, B. Denis, and X.-P. 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.

[41] L. Liu and L. Liu, "Cross-modal mutual learning for cued speech recognition," in *Proc. IEEE-ICASSP*, 2023, pp. 1–5.

[42] J. Zhang, W. Zhou, and H. Li, "A threshold-based HMM-DTW approach for continuous sign language recognition," in *Proceedings of International Conference on Internet Multimedia Computing and Service*. Association for Computing Machinery, 2014, pp. 237–240.

[43] W. Yang, J. Tao, and Z. Ye, "Continuous sign language recognition using level building based on fast hidden markov model," *Pattern Recognition Letters*, vol. 78, pp. 28–35, 2016.

[44] K. L. Cheng, Z. Yang, Q. Chen, and Y.-W. Tai, "Fully convolutional networks for continuous sign language recognition," in *Proc. ECCV*, 2020, pp. 697–714.

[45] C. Wei, J. Zhao, W. Zhou, and H. Li, "Semantic boundary detection with reinforcement learning for continuous sign language recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 3, pp. 1138–1149, 2020.

[46] L. Liu and G. Feng, "A pilot study on mandarin chinese cued speech," *American Annals of the Deaf*, vol. 164, no. 4, pp. 496–518, 2019.

[47] "Sign languages unite us!" un.org, 2022. [Online]. Available: https://www.un.org/en/observances/sign-languages-day

[48] Z. Pan, X. Qian, and H. Li, "Speaker extraction with co-speech gestures cue," *IEEE Signal Processing Letters*, vol. 29, pp. 1467–1471, 2022.

[49] C. Sondermann and M. Merkt, "Like it or learn from it: Effects of talking heads in educational videos," *Computers & Education*, vol. 193, p. 104675, 2023.

[50] W. Song, Q. He, and G. Chen, "Virtual human talking-head generation," in *Proceedings of the 2023 2nd Asia Conference on Algorithms, Computing and Machine Learning*, 2023, pp. 1–5.

[51] D. Kothadiya, C. Bhatt, K. Sapariya, K. Patel, A.-B. Gil-González, and J. M. Corchado, "Deepsign: Sign language detection and recognition using deep learning," *Electronics*, vol. 11, no. 11, p. 1780, 2022.

[52] M. De Coster, D. Shterionov, M. Van Herreweghe, and J. Dambre, "Machine translation from signed to spoken languages: State of the art and challenges," *Universal Access in the Information Society*, pp. 1–27, 2023.

[53] N. K. Kahlon and W. Singh, "Machine translation from text to sign language: a systematic review," *Universal Access in the Information Society*, vol. 22, no. 1, pp. 1–35, 2023.

[54] L. Liu, G. Feng, X. Ren, and X. Ma, "Objective hand complexity comparison between two mandarin chinese cued speech systems," in *Proc. IEEE-ISCSLP*, 2022, pp. 215–219.

[55] "Find your cued language," cuedspeech.org. [Online]. Available: https://cuedspeech.org/learn/find-your-cued-language/

[56] L. Liu, "Modeling for continuous cued speech recognition in french using advanced machine learning methods," Ph.D. dissertation, Universite Grenoble Alpes, 2018.

[57] L. Liu, G. Feng, and D. Beautemps, "Automatic temporal segmentation of hand movements for hand positions recognition in french cued speech," in *Proc. IEEE-ICASSP*, 2018, pp. 3061–3065.

[58] J. Wang, Z. Tang, X. Li, M. Yu, Q. Fang, and L. Liu, "Cross-modal knowledge distillation method for automatic cued speech recognition," in *Proc. Interspeech*, 2021, p. 2986–2990.

[59] L. Liu, J. Li, G. Feng, and X.-P. S. Zhang, "Automatic detection of the temporal segmentation of hand movements in british english cued speech." in *Proc. Interspeech*, 2019, pp. 2285–2289.

[60] S. J. Park, M. Kim, J. Hong, J. Choi, and Y. M. Ro, "Synctalkface: Talking face generation with precise lip-syncing via audio-lip memory," in *Proc. Conf AAAI Artif. Intell.*, vol. 36, no. 2, 2022, pp. 2062–2070.

[61] L. Liu, G. Feng, and D. Beautemps, "Inner lips parameter estimation based on adaptive ellipse model," in *International Conference on Auditory-Visual Speech Processing (AVSP)*, 2017.

[62] ——, "Automatic dynamic template tracking of inner lips based on clnf," in *Proc. IEEE-ICASSP*, 2017, p. 5130–5134.

[63] M. W. Alibali, M. Bassok, K. O. Solomon, S. E. Syc, and S. Goldin-Meadow, "Illuminating mental representations through speech and gesture," *Psychological Science*, vol. 10, no. 4, pp. 327–333, 1999.

[64] S. Kang and B. Tversky, "From hands to minds: Gestures promote understanding," *Cognitive Research: Principles and Implications*, vol. 1, no. 1, pp. 1–15, 2016.

[65] A. Kendon, "Do gestures communicate? a review," *Research on language and social interaction*, vol. 27, no. 3, pp. 175–200, 1994.

[66] K. Adam, *Gesture: Visible action as utterance*. Cambridge University Press, 2004.

[67] Y. Ferstl and R. McDonnell, "Investigating the use of recurrent motion modelling for speech gesture generation," in *Proc. ACM IVA*, 2018, pp. 93–98.

[68] Y. Yoon, W.-R. Ko, M. Jang, J. Lee, J. Kim, and G. Lee, "Robots learn social skills: End-to-end learning of co-speech gesture generation for humanoid robots," in *Proc. IEEE-International Conference in Robotics and Automation (ICRA)*, 2019, pp. 4303–4309.

[69] Y. Yoon, P. Wolfert, T. Kucherenko, C. Viegas, T. Nikolov, M. Tsakov, and G. E. Henter, "The genea challenge 2022: A large evaluation of data-driven co-speech gesture generation," in *Proc. ACM-International Conference on Multimodal Interaction*, 2022, pp. 736–747.

[70] S. Ginosar, A. Bar, G. Kohavi, C. Chan, A. Owens, and J. Malik, "Learning individual styles of conversational gesture," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 3497–3506.

[71] R. Poppe, "A survey on vision-based human action recognition," *Image and vision computing*, vol. 28, no. 6, pp. 976–990, 2010.

[72] D. Wu, L. Pigou, P.-J. Kindermans, N. D.-H. Le, L. Shao, J. Dambre, and J.-M. Odobez, "Deep dynamic neural networks for multimodal gesture segmentation and recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 8, pp. 1583–1597, 2016.

[73] J. Wan, Y. Zhao, S. Zhou, I. Guyon, S. Escalera, and S. Z. Li, "Chalearn looking at people rgb-d isolated and continuous datasets for gesture recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2016, pp. 56–64.

[74] J. Materzynska, G. Berger, I. Bax, and R. Memisevic, "The jester dataset: A large-scale video dataset of human gestures," in *Proceedings of the IEEE/CVF international conference on computer vision workshops*, 2019, pp. 0–0.

[75] H. Zeng, X. Wang, Y. Wang, A. Wu, T.-C. Pong, and H. Qu, "Gesturelens: Visual analysis of gestures in presentation videos," *IEEE Transactions on Visualization and Computer Graphics*, 2022.

[76] E. Efthimiou, S.-E. Fotinea, T. Hanke, J. Glauert, R. Bowden, A. Braffort, C. Collet, P. Maragos, and F. Goudenove, "Dicta-sign: sign language recognition, generation and modelling with application in deaf communication," in *LREC*. European Language Resources Association (ELRA), 2010, pp. 80–83.

[77] O. Koller, J. Forster, and H. Ney, "Continuous sign language recognition: Towards large vocabulary statistical recognition systems handling multiple signers," *Computer Vision and Image Understanding*, vol. 141, pp. 108–125, 2015.

[78] C. Neidle, A. Thangali, and S. Sclaroff, "Challenges in development of the american sign language lexicon video dataset (asllvd) corpus," in *LREC*. Citeseer, 2012.

[79] U. v. Agris and K.-F. Kraiss, "Signum database: Video corpus for signer-independent continuous sign language recognition," in *LREC*. European Language Resources Association (ELRA), 2010, pp. 243–246.

[80] Y. Lin, X. Chai, Y. Zhou, and X. Chen, "Curve matching from the view of manifold for sign language recognition," in *ACCV Workshops*, 2014.

[81] N. K. Caselli, Z. S. Sehyr, A. M. Cohen-Goldberg, and K. Emmorey, "Asl-lex: A lexical database of american sign language," *Behavior research methods*, vol. 49, pp. 784–801, 2017.

[82] N. C. Camgoz, S. Hadfield, O. Koller, H. Ney, and R. Bowden, "Neural sign language translation," in *Proc. IEEE-CVPR*, 2018, pp. 7784–7793.

[83] A. Mavi and Z. Dikle, "A new 27 class sign language dataset collected from 173 individuals," *arXiv preprint arXiv:2203.03859*, 2022.

[84] S.-K. Ko, C. J. Kim, H. Jung, and C. Cho, "Neural sign language translation based on human keypoint estimation," *Applied sciences*, vol. 9, no. 13, p. 2683, 2019.

[85] N. Adaloglou, T. Chatzis, I. Papastratis, A. Stergioulas, G. T. Papadopoulos, V. Zacharopoulou, G. J. Xydopoulos, K. Atzakas, D. Papazachariou, and P. Daras, "A comprehensive study on deep learning-based methods for sign language recognition," *IEEE Transactions on Multimedia*, vol. 24, pp. 1750–1762, 2021.

[86] Z. S. Sehyr, N. Caselli, A. M. Cohen-Goldberg, and K. Emmorey, "The asl-lex 2.0 project: A database of lexical and phonological properties for 2,723 signs in american sign language," *The Journal of Deaf Studies and Deaf Education*, vol. 26, no. 2, pp. 263–277, 2021.

[87] A. Duarte, S. Palaskar, L. Ventura, D. Ghadiyaram, K. DeHaan, F. Metze, J. Torres, and X. Giro-i Nieto, "How2sign: a large-scale multimodal dataset for continuous american sign language," in *Proc. IEEE/CVF-CVPR*, 2021, pp. 2735–2744.

[88] A. Kapitanov, K. Kvanchiani, A. Nagaev, and E. Petrova, "Slovo: Russian sign language dataset," *arXiv preprint arXiv:2305.14527*, 2023.

[89] M. Al-Barham, A. Alsharkawi, M. Al-Yaman, M. Al-Fetyani, A. El-nagar, A. A. SaAleek, and M. Al-Odat, "Rgb arabic alphabets sign language dataset," *arXiv preprint arXiv:2301.11932*, 2023.

[90] L. Liu, H. Thomas, G. Feng, and B. Denis, "Visual recognition of continuous cued speech using a tandem cnn-hmm approach." in *Proc. Interspeech*, 2018, pp. 2643–2647.

[91] T. A., "VOT and durational properties of selected segments in the speech of deaf and normally hearing children," *Studia Phonetica Posnaniensia*, vol. 8, pp. 111–142, 2007.

[92] B. Bigi, M. Zimmermann, and C. André, "Clelfpc: a large open multi-speaker corpus of french cued speech," in *LREC*, 2022, pp. 987–994.

[93] M. Cooke, J. Barker, S. Cunningham, and X. Shao, "The grid audio-visual speech corpus (1.0) [data set]," in *Zenodo*. Zenodo, 2006.

[94] O. Martin, I. Kotsia, B. Macq, and I. Pitas, "The enterface'05 audio-visual emotion database," in *Proc. IEEE-22nd international conference on data engineering workshops*, 2006, pp. 8–8.

[95] A. Rekik, A. Ben-Hamadou, and W. Mahdi, "An adaptive approach for lip-reading using image and depth data," *Multimedia Tools and Applications*, vol. 75, pp. 8609–8636, 2016.

[96] H. Cao, D. G. Cooper, M. K. Keutmann, R. C. Gur, A. Nenkova, and R. Verma, "Crema-d: Crowd-sourced emotional multimodal actors dataset," *IEEE Transactions on Affective Computing*, vol. 5, no. 4, pp. 377–390, 2014.

[97] A. Czyzewski, B. Kostek, P. Bratoszewski, J. Kotus, and M. Szykulski, "An audio-visual corpus for multimodal automatic speech recognition," *Journal of Intelligent Information Systems*, vol. 49, pp. 167–192, 2017.

[98] J. S. Chung and A. Zisserman, "Lip reading in the wild," *Proc. ACCV*, pp. 87–103, 2017.

[99] C. Busso, S. Parthasarathy, A. Burmania, M. AbdelWahab, N. Sadoughi, and E. M. Provost, "MSP-IMPROV: An acted corpus of dyadic interactions to study emotion perception," *IEEE Transactions on Affective Computing*, vol. 8, no. 1, pp. 67–80, 2016.

[100] S. Suwajanakorn, S. M. Seitz, and I. Kemelmacher-Shlizerman, "Synthesizing obama: learning lip sync from audio," *ACM Transactions on Graphics (ToG)*, vol. 36, no. 4, pp. 1–13, 2017.

[101] A. Nagrani, J. S. Chung, and A. Zisserman, "Voxceleb: a large-scale speaker identification dataset," *Telephony*, vol. 3, pp. 33–039, 2017.

[102] J. S. Chung, A. Nagrani, and A. Zisserman, "Voxceleb2: Deep speaker recognition," *Proc. Interspeech*, 2018.

[103] 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.

[104] T. Afouras, J. S. Chung, and A. Zisserman, "LRS3-TED: a large-scale dataset for visual speech recognition," *arXiv preprint arXiv:1809.00496*, 2018.

[105] S. R. Livingstone and F. A. Russo, "The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english," *PloS one*, vol. 13, no. 5, p. e0196391, 2018.

[106] S. Poria, D. Hazarika, N. Majumder, G. Naik, E. Cambria, and R. Mihalcea, "Meld: A multimodal multi-party dataset for emotion recognition in conversations," *arXiv preprint arXiv:1810.02508*, 2018.

[107] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, "Looking to listen at the cocktail party: a speaker-independent audio-visual model for speech separation," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 4, pp. 1–11, 2018.

[108] D. Cudeiro, T. Bolkart, C. Laidlaw, A. Ranjan, and M. J. Black, "Capture, learning, and synthesis of 3d speaking styles," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 10 101–10 111.

[109] S. Yang, Y. Zhang, D. Feng, M. Yang, C. Wang, J. Xiao, K. Long, S. Shan, and X. Chen, "LRW-1000: A naturally-distributed large-scale benchmark for lip reading in the wild," in *Proceedings of 14th IEEE international conference on automatic face & gesture recognition*, 2019, pp. 1–8.

[110] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "Faceforensics++: Learning to detect manipulated facial images," in *Proc. IEEE/CVF-ICCV*, 2019, pp. 1–11.

[111] K. Wang, Q. Wu, L. Song, Z. Yang, W. Wu, C. Qian, R. He, Y. Qiao, and C. C. Loy, "Mead: A large-scale audio-visual dataset for emotional talking-face generation," in *Proc. ECCV*, 2020, pp. 700–717.

[112] Z. Zhang, L. Li, Y. Ding, and C. Fan, "Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset," in *Proc. IEEE/CVF-CVPR*, 2021, pp. 3661–3670.

[113] K. Kim, S. Park, J. Lee, S. Chung, J. Lee, and J. Choo, "AnimeCeleb: Large-scale animation celebheads dataset for head reenactment," in *Proc. ECCV*, 2022, pp. 414–430.

[114] A. Berkol, T. Tümer-Sivri, N. Pervan-Akman, M. Çolak, and H. Erdem, "Visual lip reading dataset in turkish," *Data*, vol. 8, no. 1, p. 15, 2023.

[115] G. Hwang, S. Hong, S. Lee, S. Park, and G. Chae, "DisCoHead: Audio-and-video-driven talking head generation by disentangled control of head pose and facial expressions," in *Proc. IEEE-ICASSP*, 2023, pp. 1–5.

[116] I. Matthews, T. F. Cootes, J. A. Bangham, S. Cox, and R. Harvey, "Extraction of visual features for lipreading," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 2, pp. 198–213, 2002.

[117] S. Petridis, J. Shen, D. Cetin, and M. Pantic, "Visual-only recognition of normal, whispered and silent speech," in *Proc. IEEE-ICASSP*, 2018, pp. 6219–6223.

[118] K. Takeuchi, S. Kubota, K. Suzuki, D. Hasegawa, and H. Sakuta, "Creating a gesture-speech dataset for speech-based automatic gesture generation," *Communications in Computer and Information Science*, pp. 198–202, 2017.

[119] N. Singh, J. J. Lee, I. Grover, and C. Breazeal, "P2pstory: dataset of children as storytellers and listeners in peer-to-peer interactions," in *Proceedings of the CHI Conference on Human Factors in Computing Systems*, 2018, pp. 1–11.

[120] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black, "AMASS: Archive of motion capture as surface shapes," in *Proc. IEEE/CVF-ICCV*, 2019, pp. 5442–5451.

[121] Y. Luo, J. Ye, R. B. Adams, J. Li, M. G. Newman, and J. Z. Wang, "ARBEE: Towards automated recognition of bodily expression of emotion in the wild," *International journal of computer vision*, vol. 128, pp. 1–25, 2020.

[122] C. Ahuja, D. W. Lee, R. Ishii, and L.-P. Morency, "No gestures left behind: Learning relationships between spoken language and freeform gestures," in *Findings of the Association for Computational Linguistics: EMNLP*, 2020, pp. 1884–1895.

[123] A. R. Punnakkal, A. Chandrasekaran, N. Athanasiou, A. Quiros-Ramirez, and M. J. Black, "Babel: Bodies, action and behavior with english labels," in *Proc. IEEE/CVF-CVPR*, 2021, pp. 722–731.

[124] C. Guo, S. Zou, X. Zuo, S. Wang, W. Ji, X. Li, and L. Cheng, "Generating diverse and natural 3d human motions from text," in *Proc. IEEE/CVF-CVPR*, 2022, pp. 5152–5161.

[125] H. Liu, Z. Zhu, N. Iwamoto, Y. Peng, Z. Li, Y. Zhou, E. Bozkurt, and B. Zheng, "Beat: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis," in *Proc. ECCV*, 2022, pp. 612–630.

[126] J. Wang, Y. Zhao, L. Liu, T. Xu, Q. Li, and S. Li, "Emotional talking head generation based on memory-sharing and attention-augmented networks," *arXiv preprint arXiv:2306.03594*, 2023.

[127] J. Wang, Y. Zhao, H. Fan, T. Xu, Q. Li, S. Li, and L. Liu, "Memory-augmented contrastive learning for talking head generation," in *Proc. IEEE-ICASSP*, 2023, p. 1–5.

[128] E. Cosatto, J. Ostermann, H. P. Graf, and J. Schroeter, "Lifelike talking faces for interactive services," *Proceedings of the IEEE*, vol. 91, no. 9, pp. 1406–1429, 2003.

[129] O. Gambino, A. Augello, A. Caronia, G. Pilato, R. Pirrone, and S. Gaglio, "Virtual conversation with a real talking head," in *Proc. IEEE-Conference on Human System Interactions*, 2008, pp. 263–268.

[130] Y. Yang and D. Ramanan, "Articulated human detection with flexible mixtures of parts," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 12, pp. 2878–2890, 2012.

[131] Y. Yoon, B. Cha, J.-H. Lee, M. Jang, J. Lee, J. Kim, and G. Lee, "Speech gesture generation from the trimodal context of text, audio, and speaker identity," *ACM Transactions on Graphics (TOG)*, vol. 39, no. 6, pp. 1–16, 2020.

[132] E. Asakawa, N. Kaneko, D. Hasegawa, and S. Shirakawa, "Evaluation of text-to-gesture generation model using convolutional neural network," *Neural Networks*, vol. 151, pp. 365–375, 2022.

[133] P. Buehler, A. Zisserman, and M. Everingham, "Learning sign language by watching tv (using weakly aligned subtitles)," in *Proc. IEEE-CVPR*, 2009, pp. 2961–2968.

[134] H. Wang, X. Chai, and X. Chen, "A novel sign language recognition framework using hierarchical grassmann covariance matrix," *IEEE Transactions on Multimedia*, vol. 21, no. 11, pp. 2806–2814, 2019.

[135] T. Pfister, J. Charles, and A. Zisserman, "Large-scale learning of sign language by watching tv (using co-occurrences)." in *Proc. BMVC*. British Machine Vision Association, 2013.

[136] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE-CVPR*, 2016, pp. 770–778.

[137] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *Proc. IEEE-CVPR*, 2017, pp. 6299–6308.

[138] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3d residual networks," in *Proc. IEEE-ICCV*, 2017, pp. 5533–5541.

[139] Z. Qiu, T. Yao, C.-W. Ngo, X. Tian, and T. Mei, "Learning spatio-temporal representation with local and global diffusion," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 12 056–12 065.

[140] H. Hu, J. Pu, W. Zhou, and H. Li, "Collaborative multilingual continuous sign language recognition: A unified framework," *IEEE Transactions on Multimedia*, 2022.

[141] R. Cui, H. Liu, and C. Zhang, "A deep neural framework for continuous sign language recognition by iterative training," *IEEE Transactions on Multimedia*, vol. 21, no. 7, pp. 1880–1891, 2019.

[142] J. Pu, W. Zhou, H. Hu, and H. Li, "Boosting continuous sign language recognition via cross modality augmentation," in *Proc. ACM MM*, 2020, pp. 1497–1505.

[143] J. Pu, W. Zhou, and H. Li, "Iterative alignment network for continuous sign language recognition," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 4165–4174.

[144] J. Huang, W. Zhou, Q. Zhang, H. Li, and W. Li, "Video-based sign language recognition without temporal segmentation," in *Proc. Conf AAAI Artif. Intell.*, vol. 32, no. 1, 2018.

[145] H. Hu, W. Zhou, J. Pu, and H. Li, "Global-local enhancement network for nmf-aware sign language recognition," *ACM transactions on multimedia computing, communications, and applications (TOMM)*, vol. 17, no. 3, pp. 1–19, 2021.

[146] C. Wei, W. Zhou, J. Pu, and H. Li, "Deep grammatical multi-classifier for continuous sign language recognition," in *International Conference on Multimedia Big Data (BigMM)*, 2019, pp. 435–442.

[147] D. Guo, S. Wang, Q. Tian, and M. Wang, "Dense temporal convolution network for sign language translation." in *Proc.IJCAI*, 2019, pp. 744–750.

[148] H. Zhou, W. Zhou, and H. Li, "Dynamic pseudo label decoding for continuous sign language recognition," in *International conference on multimedia and expo (ICME)*, 2019, pp. 1282–1287.

[149] S. NadeemHashmi, H. Gupta, D. Mittal, K. Kumar, A. Nanda, and S. Gupta, "A lip reading model using cnn with batch normalization," in *Proc. IEEE-11th international conference on contemporary computing (IC3)*, 2018, pp. 1–6.

[150] F. B. Slimane and M. Bouguessa, "Context matters: Self-attention for sign language recognition," in *International Conference on Pattern Recognition (ICPR)*, 2021, pp. 7884–7891.

[151] M. Zhou, M. Ng, Z. Cai, and K. C. Cheung, "Self-attention-based fully-inception networks for continuous sign language recognition," in *24th European Conference on Artificial Intelligence*, 2020, pp. 2832–2839.

[152] O. Koller, N. C. Camgoz, H. Ney, and R. Bowden, "Weakly supervised learning with multi-stream cnn-lstm-hmms to discover sequential parallelism in sign language videos," *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 9, pp. 2306–2320, 2019.

[153] O. Koller, S. Zargaran, and H. Ney, "Re-sign: Re-aligned end-to-end sequence modelling with deep recurrent CNN-HMMs," in *Proc. IEEE-CVPR*, 2017, pp. 4297–4305.

[154] O. Koller, S. Zargaran, H. Ney, and R. Bowden, "Deep sign: Enabling robust statistical continuous sign language recognition via hybrid cnn-hmms," *International Journal of Computer Vision*, vol. 126, pp. 1311–1325, 2018.

[155] R. Cui, H. Liu, and C. Zhang, "Recurrent convolutional neural networks for continuous sign language recognition by staged optimization," in *Proc. IEEE-CVPR*, 2017, pp. 7361–7369.

[156] Y. Min, A. Hao, X. Chai, and X. Chen, "Visual alignment constraint for continuous sign language recognition," in *Proc. IEEE/CVF-ICCV*, 2021, pp. 11 542–11 551.

[157] N. Cihan Camgoz, S. Hadfield, O. Koller, and R. Bowden, "Subunets: End-to-end hand shape and continuous sign language recognition," in *Proc. IEEE-ICCV*, 2017, pp. 3056–3065.

[158] D. Guo, W. Zhou, A. Li, H. Li, and M. Wang, "Hierarchical recurrent deep fusion using adaptive clip summarization for sign language translation," *IEEE Transactions on Image Processing*, vol. 29, pp. 1575–1590, 2019.

[159] H. Li, L. Gao, R. Han, L. Wan, and W. Feng, "Key action and joint ctc-attention based sign language recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 2348–2352.

[160] A. Hao, Y. Min, and X. Chen, "Self-mutual distillation learning for continuous sign language recognition," in *Proc. IEEE/CVF-ICCV*, 2021, pp. 11 303–11 312.

[161] H. Zhou, W. Zhou, Y. Zhou, and H. Li, "Spatial-temporal multi-cue network for sign language recognition and translation," *IEEE Transactions on Multimedia*, vol. 24, pp. 768–779, 2021.

[162] X. Pei, D. Guo, and Y. Zhao, "Continuous sign language recognition based on pseudo-supervised learning," in *Proceedings of the 2nd Workshop on Multimedia for Accessible Human Computer Interfaces*, 2019, pp. 33–39.

[163] Z. Zhang, J. Pu, L. Zhuang, W. Zhou, and H. Li, "Continuous sign language recognition via reinforcement learning," in *Proc. IEEE-ICIP*. IEEE, 2019, pp. 285–289.

[164] Z. Niu and B. Mak, "Stochastic fine-grained labeling of multi-state sign glosses for continuous sign language recognition," in *Proc. ECCV*, 2020, pp. 172–186.

[165] K. Koishybay, M. Mukushev, and A. Sandygulova, "Continuous sign language recognition with iterative spatiotemporal fine-tuning," in *Proc. IEEE-ICPR*, 2021, pp. 10 211–10 218.

[166] I. Papastratis, K. Dimitropoulos, and P. Daras, "Continuous sign language recognition through a context-aware generative adversarial network," *Sensors*, vol. 21, no. 7, 2021.

[167] Y. Chen, R. Zuo, F. Wei, Y. Wu, S. Liu, and B. Mak, "Two-stream network for sign language recognition and translation," *Advances in Neural Information Processing Systems (NIPS)*, vol. 35, pp. 17 043–17 056, 2022.

[168] L. Hu, L. Gao, Z. Liu, and W. Feng, "Self-emphasizing network for continuous sign language recognition," in *Proc. Conf AAAI Artif. Intell.*, vol. 37, no. 1, 2023, pp. 854–862.

[169] J. Zheng, Y. Wang, C. Tan, S. Li, G. Wang, J. Xia, Y. Chen, and S. Z. Li, "Cvt-slr: Contrastive visual-textual transformation for sign language recognition with variational alignment," in *Proc. IEEE/CVF-CVPR*, 2023, pp. 23 141–23 150.

[170] K. Papadimitriou and G. Potamianos, "A fully convolutional sequence learning approach for cued speech recognition from videos," in *Proc. IEEE-EUSIPCO*, 2021, pp. 326–330.

[171] S. Sankar, D. Beautemps, and T. Hueber, "Multistream neural architectures for cued speech recognition using a pre-trained visual feature extractor and constrained ctc decoding," in *Proc. IEEE-ICASSP*, 2022, pp. 8477–8481.

[172] M. Kipp, A. Heloir, and Q. Nguyen, "Sign language avatars: Animation and comprehensibility," in *Intelligent Virtual Agents*. Springer Berlin Heidelberg, 2011, pp. 113–126.

[173] J. McDonald, R. Wolfe, J. Schnepp, J. Hochgesang, D. G. Jamrozik, M. Stumbo, L. Berke, M. Bialek, and F. Thomas, "An automated technique for real-time production of lifelike animations of american sign language," *Universal Access in the Information Society*, vol. 15, pp. 551–566, 2016.

[174] S. Gibet, F. Lefebvre-Albaret, L. Hamon, R. Brun, and A. Turki, "Interactive editing in french sign language dedicated to virtual signers: Requirements and challenges," *Universal Access in the Information Society*, vol. 15, pp. 525–539, 2016.

[175] J. Zelinka and J. Kanis, "Neural sign language synthesis: Words are our glosses," in *Proc. IEEE/CVF-WACV*, March 2020.

[176] N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, "Multi-channel transformers for multi-articulatory sign language translation," 2020.

[177] B. Saunders, N. C. Camgöz, and R. Bowden, "Adversarial training for multi-channel sign language production," in *The 31st British Machine Vision Virtual Conference*. British Machine Vision Association, 2020.

[178] M. Inan, Y. Zhong, S. Hassan, L. Quandt, and M. Alikhani, "Modeling intensification for sign language generation: A computational approach," in *Findings of the Association for Computational Linguistics*, 2022, pp. 2897–2911.

[179] B. Saunders, N. C. Camgoz, and R. Bowden, "Signing at scale: Learning to co-articulate signs for large-scale photo-realistic sign language production," in *Proc. IEEE/CVF-CVPR*, 2022, pp. 5141–5151.

[180] P. Xie, Q. Zhang, Z. Li, H. Tang, Y. Du, and X. Hu, "Vector quantized diffusion model with codeunet for text-to-sign pose sequences generation," *arXiv preprint arXiv:2208.09141*, 2022.

[181] C.-C. Chiu, L.-P. Morency, and S. Marsella, "Predicting co-verbal gestures: A deep and temporal modeling approach," in *Proc. Intelligent Virtual Agents*. Springer International Publishing, 2015, pp. 152–166.

[182] S. Alexanderson, G. E. Henter, T. Kucherenko, and J. Beskow, "Style-controllable speech-driven gesture synthesis using normalising flows," in *Computer Graphics Forum*, vol. 39, no. 2. Wiley Online Library, 2020, pp. 487–496.

[183] J. Li, D. Kang, W. Pei, X. Zhe, Y. Zhang, Z. He, and L. Bao, "Audio2gestures: Generating diverse gestures from speech audio with conditional variational autoencoders," in *Proc. IEEE/CVF-ICCV*, 2021, pp. 11 293–11 302.

[184] U. Bhattacharya, N. Rewkowski, A. Banerjee, P. Guhan, A. Bera, and D. Manocha, "Text2gestures: A transformer-based network for generating emotive body gestures for virtual agents," in *IEEE virtual reality and 3D user interfaces (VR)*, 2021, pp. 1–10.

[185] S. Ghorbani, Y. Ferstl, and M.-A. Carbonneau, "Exemplar-based stylized gesture generation from speech: An entry to the GENEA challenge 2022," in *Proc. ACM-International Conference on Multimodal Interaction*, 2022, pp. 778–783.

[186] M. Li and Y.-m. Cheung, "A novel motion based lip feature extraction for lip-reading," in *Proc. IEEE-International Conference on Computational Intelligence and Security*, vol. 1, 2008, pp. 361–365.

[187] S. Alizadeh, R. Boostani, and V. Asadpour, "Lip feature extraction and reduction for hmm-based visual speech recognition systems," in *Pro. IEEE-9th International Conference on Signal Processing*, 2008, pp. 561–564.

[188] X. Ma, L. Yan, and Q. Zhong, "Lip feature extraction based on improved jumping-snake model," in *Proc. IEEE-35th Chinese Control Conference (CCC)*, 2016, pp. 6928–6933.

[189] Y. Lan, B.-J. Theobald, and R. Harvey, "View independent computer lip-reading," in *Proc. IEEE-International Conference on Multimedia and Expo*, 2012, pp. 432–437.

[190] T. Watanabe, K. Katsurada, and Y. Kanazawa, "Lip reading from multi view facial images using 3D-AAM," in *Proc. ACCV*. Springer Verlag, 2017, pp. 303–316.

[191] L. Liu, G. Feng, and D. Beautemps, "Inner lips feature extraction based on clnf with hybrid dynamic template for cued speech," *EURASIP Journal on Image and Video Processing*, vol. 2017, p. 1–15, 2017.

[192] ——, "Extraction automatique de contour de levre a partir du modele clnf," in *JEP-TALN-RECITAL 2016-conference conjointe 31e Journees d'Etudes sur la Parole, 23e Traitement Automatique des Langues Naturelles, 18e Rencontre des Etudiants Chercheurs en Informatique pour le Traitement Automatique des Langues*, 2016.

[193] J. Wang, T. Wu, S. Wang, M. Yu, Q. Fang, J. Zhang, and L. Liu, "Three-dimensional lip motion network for text-independent speaker recognition," in *Proc. IEEE-ICPR*, 2021, p. 3380–3387.

[194] A. Garg, J. Noyola, and S. Bagadia, "Lip reading using cnn and lstm," *Technical report, Stanford University, CS231 n project report*, 2016.

[195] D. Lee, J. Lee, and K.-E. Kim, "Multi-view automatic lip-reading using neural network," in *Proc. ACCV Workshop on Multi-view Lip-reading Challenges*, 2016.

[196] I. Fung and B. Mak, "End-to-end low-resource lip-reading with maxout CNN and LSTM," in *Proc. IEEE-ICASSP*, 2018, pp. 2511–2515.

[197] K. Xu, D. Li, N. Cassimatis, and X. Wang, "LCANet: End-to-end lipreading with cascaded attention-CTC," in *Proc. IEEE-FG*, 2018, pp. 548–555.

[198] P. Wiriyathammabhum, "Spotfast networks with memory augmented lateral transformers for lipreading," in *International Conference on Neural Information Processing*, 2020, pp. 554–561.

[199] X. Weng and K. Kitani, "Learning spatio-temporal features with two-stream deep 3d cnns for lipreading," *arXiv preprint arXiv:1905.02540*, 2019.

[200] D. Feng, S. Yang, and S. Shan, "An efficient software for building lip reading models without pains," in *Proc. IEEE-ICMEW*, 2021, pp. 1–2.

[201] B. Xu, C. Lu, Y. Guo, and J. Wang, "Discriminative multi-modality speech recognition," in *Proc. IEEE/CVF-CVPR*, 2020, pp. 14 433–14 442.

[202] M. Luo, S. Yang, S. Shan, and X. Chen, "Pseudo-convolutional policy gradient for sequence-to-sequence lip-reading," in *Proc. IEEE-FG*, 2020, pp. 273–280.

[203] J. Gehring, Y. Miao, F. Metze, and A. Waibel, "Extracting deep bottleneck features using stacked auto-encoders," in *Proc. IEEE-ICASSP*, 2013, pp. 3377–3381.

[204] K. Noda, Y. Yamaguchi, K. Nakadai, H. G. Okuno, and T. Ogata, "Audio-visual speech recognition using deep learning," *Applied intelligence*, vol. 42, pp. 722–737, 2015.

[205] S. Petridis and M. Pantic, "Deep complementary bottleneck features for visual speech recognition," in *Proc. IEEE-ICASSP*, 2016, pp. 2304–2308.

[206] J. Wang, N. Gu, M. Yu, X. Li, Q. Fang, and L. Liu, "An attention self-supervised contrastive learning based three-stage model for hand shape feature representation in cued speech," *arXiv preprint arXiv:2106.14016*, 2021.

[207] M. Wand, J. Schmidhuber, and N. T. Vu, "Investigations on end-to-end audiovisual fusion," in *Proc. IEEE-ICASSP*, 2018, pp. 3041–3045.

[208] Y. Zhang, S. Yang, J. Xiao, S. Shan, and X. Chen, "Can we read speech beyond the lips? rethinking roi selection for deep visual speech recognition," in *Proc. IEEE-FG*, 2020, pp. 356–363.

[209] J. Xiao, S. Yang, Y. Zhang, S. Shan, and X. Chen, "Deformation flow based two-stream network for lip reading," in *Proc. IEEE-FG*, 2020, pp. 364–370.

[210] X. Zhao, S. Yang, S. Shan, and X. Chen, "Mutual information maximization for effective lip reading," in *Proc. IEEE-FG*, 2020, pp. 420–427.

[211] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.

[212] B. Martinez, P. Ma, S. Petridis, and M. Pantic, "Lipreading using temporal convolutional networks," in *Proc. IEEE-ICASSP*, 2020, pp. 6319–6323.

[213] V. Ashish, S. Noam, P. Niki, U. Jakob, G. A. N, K. Łukasz, and P. Illia, "Attention is all you need," in *Advances in Neural Information Processing Systems (NIPS)*, 2017.

[214] T. Afouras, J. S. Chung, and A. Zisserman, "Deep lip reading: a comparison of models and an online application," *arXiv preprint arXiv:1806.06053*, 2018.

[215] J. Son Chung, A. Senior, O. Vinyals, and A. Zisserman, "Lip reading sentences in the wild," in *Proc. IEEE-CVPR*, 2017, pp. 6447–6456.

[216] Y. Lu and H. Li, "Automatic lip-reading system based on deep convolutional neural network and attention-based long short-term memory," *Applied Sciences*, vol. 9, no. 8, p. 1599, 2019.

[217] P. Zhou, W. Yang, W. Chen, Y. Wang, and J. Jia, "Modality attention for end-to-end audio-visual speech recognition," in *Proc. IEEE-ICASSP*, 2019, pp. 6565–6569.

[218] X. Zhang, H. Gong, X. Dai, F. Yang, N. Liu, and M. Liu, "Understanding pictograph with facial features: End-to-end sentence-level lip reading of chinese," in *Proc. Conf AAAI Artif. Intell.*, vol. 33, no. 01, 2019, pp. 9211–9218.

[219] A. Torfi, S. M. Iranmanesh, N. Nasrabadi, and J. Dawson, "3d convolutional neural networks for cross audio-visual matching recognition," *IEEE Access*, vol. 5, pp. 22 081–22 091, 2017.

[220] P. Heracleous, D. Beautemps, and N. Aboutabit, "Cued speech automatic recognition in normal-hearing and deaf subjects," *Speech Communication*, vol. 52, no. 6, pp. 504–512, 2010.

[221] P. Heracleous, D. Beautemps, and N. Hagita, "Continuous phoneme recognition in cued speech for french," in *Proc. IEEE-EUSIPCO*, 2012, pp. 2090–2093.

[222] T. Burger, A. Caplier, and S. Mancini, "Cued speech hand gestures recognition tool," in *Proc. IEEE-EUSIPCO*, 2005, pp. 1–4.

[223] L. Gao, S. Huang, and L. Liu, "A novel interpretable and generalizable re-synchronization model for cued speech based on a multi-cuer corpus," *arXiv preprint arXiv:2306.02596*, 2023.

[224] Y. Zhang, L. Liu, and L. Liu, "Cuing without sharing: A federated cued speech recognition framework via mutual knowledge distillation," *arXiv preprint arXiv:2308.03432*, 2023.

[225] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, 2006, pp. 369–376.

[226] R. Boeck, K. Bergmann, and P. Jaecks, "Disposition recognition from spontaneous speech towards a combination with co-speech gestures," in *Proceedings of the 2nd International Workshop on Multimodal Analyses enabling Artificial Agents in Human-Machine Interaction*, 2014.

[227] U. Bhattacharya, E. Childs, N. Rewkowski, and D. Manocha, "Speech2affectivegestures: Synthesizing co-speech gestures with generative adversarial affective expression learning," in *Proc. ACM MM*, 2021, pp. 2027–2036.

[228] Y. Wen, B. Raj, and R. Singh, "Face reconstruction from voice using generative adversarial networks," *Advances in Neural Information Processing Systems (NIPS)*, vol. 32, 2019.

[229] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proc. IEEE-CVPR*, 2015, pp. 815–823.

[230] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 4690–4699.

[231] B. Saunders, N. C. Camgoz, and R. Bowden, "Progressive transformers for end-to-end sign language production," in *Proc. ECCV*, 2020, pp. 687–705.

[232] S. Stoll, N. C. Camgöz, S. Hadfield, and R. Bowden, "Sign language production using neural machine translation and generative adversarial networks," in *Proc. BMVC*. British Machine Vision Association, 2018.

[233] N. Vasani, P. Autee, S. Kalyani, and R. Karani, "Generation of indian sign language by sentence processing and generative adversarial networks," in *Proc. IEEE-ICISS*, 2020, pp. 1250–1255.

[234] L. Ventura, A. Duarte, and X. Giró-i Nieto, "Can everybody sign now? exploring sign language video generation from 2d poses," *arXiv preprint arXiv:2012.10941*, 2020.

[235] Q. Xiao, M. Qin, and Y. Yin, "Skeleton-based chinese sign language recognition and generation for bidirectional communication between deaf and hearing people," *Neural networks*, vol. 125, pp. 41–55, 2020.

[236] E. Rothauser, "Ieee recommended practice for speech quality measurements," *IEEE Transactions on Audio and Electroacoustics*, vol. 17, no. 3, pp. 225–246, 1969.

[237] J. Kim, J. Kim, and S. Choi, "Flame: Free-form language-based motion synthesis & editing," in *Proc. Conf AAAI Artif. Intell.*, vol. 37, no. 7, 2023, pp. 8255–8263.

[238] J. Cassell, H. H. Vilhjálmsson, and T. Bickmore, "Beat: the behavior expression animation toolkit," in *Proc. ACM SIGGRAPH*, 2001, pp. 477–486.

[239] J. Cassell, C. Pelachaud, N. Badler, M. Steedman, B. Achorn, T. Becket, B. Douville, S. Prevost, and M. Stone, "Animated conversation: rule-based generation of facial expression, gesture & spoken intonation for multiple conversational agents," in *Proc. ACM SIGGRAPH*, 1994, pp. 413–420.

[240] P. Wagner, Z. Malisz, and S. Kopp, "Gesture and speech in interaction: An overview," *Speech Communication*, vol. 57, pp. 209–232, 2014.

[241] S. Levine, P. Krähenbühl, S. Thrun, and V. Koltun, "Gesture controllers," *ACM Transactions on Graphics (TOG)*, vol. 29, no. 4, pp. 1–11, 2010.

[242] S. Qian, Z. Tu, Y. Zhi, W. Liu, and S. Gao, "Speech drives templates: Co-speech gesture synthesis with learned templates," in *Proc. IEEE/CVF-ICCV*, 2021, pp. 11 077–11 086.

[243] P. Dhariwal and A. Q. Nichol, "Diffusion models beat GANs on image synthesis," in *Advances in Neural Information Processing Systems (NIPS)*, 2021.

[244] L. Zhu, X. Liu, X. Liu, R. Qian, Z. Liu, and L. Yu, "Taming diffusion models for audio-driven co-speech gesture generation," in *Proc. IEEE/CVF-CVPR*, 2023, pp. 10 544–10 553.

[245] H. Teager and S. Teager, "Evidence for nonlinear sound production mechanisms in the vocal tract," *Speech production and speech modelling*, pp. 241–261, 1990.

[246] T.-H. Oh, T. Dekel, C. Kim, I. Mosseri, W. T. Freeman, M. Rubinstein, and W. Matusik, "Speech2face: Learning the face behind a voice," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 7539–7548.

[247] O. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," in *Proc. BMVC*. British Machine Vision Association, 2015.

[248] A. C. Duarte, F. Roldan, M. Tubau, J. Escur, S. Pascual, A. Salvador, E. Mohedano, K. McGuinness, J. Torres, and X. Giro-i Nieto, "Wav2pix: Speech-conditioned face generation using generative adversarial networks." in *Proc. IEEE-ICASSP*, 2019, pp. 8633–8637.

[249] S. Pascual, A. Bonafonte, and J. Serrà, "SEGAN: Speech enhancement generative adversarial network," *Proc. Interspeech*, pp. 3642–3646, 2017.

[250] Z. Fang, Z. Liu, T. Liu, C.-C. Hung, J. Xiao, and G. Feng, "Facial expression gan for voice-driven face generation," *The Visual Computer*, pp. 1–14, 2022.

[251] F. Cole, D. Belanger, D. Krishnan, A. Sarna, I. Mosseri, and W. T. Freeman, "Synthesizing normalized faces from facial identity features," in *Proc. IEEE-CVPR*, 2017, pp. 3703–3712.

[252] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.

[253] J. Wang, Z. Wang, X. Hu, X. Li, Q. Fang, and L. Liu, "Residual-guided personalized speech synthesis based on face image," in *Proc. IEEE-ICASSP*, 2022, p. 4743–4747.

[254] D. E. King, "Dlib-ml: A machine learning toolkit," *The Journal of Machine Learning Research*, vol. 10, pp. 1755–1758, 2009.

[255] L. Chen, Z. Li, R. K. Maddox, Z. Duan, and C. Xu, "Lip movements generation at a glance," in *Proc. ECCV*, 2018, pp. 520–535.

[256] Y. Song, J. Zhu, D. Li, A. Wang, and H. Qi, "Talking face generation by conditional recurrent adversarial network," in *Proc. IJCAI*, 2019, pp. 919–925.

[257] H. Zhou, Y. Liu, Z. Liu, P. Luo, and X. Wang, "Talking face generation by adversarially disentangled audio-visual representation," in *Proc. Conf AAAI Artif. Intell.*, vol. 33, no. 01, 2019, 9299–9306.

[258] L. Chen, R. K. Maddox, Z. Duan, and C. Xu, "Hierarchical cross-modal talking face generation with dynamic pixel-wise loss," in *Proc. IEEE/CVF-CVPR*, 2019, pp. 7832–7841.

[259] K. Vougioukas, S. Petridis, and M. Pantic, "End-to-end speech-driven realistic facial animation with temporal GANs." in *CVPR Workshops*, 2019, pp. 37–40.

[260] T. Kefalas, K. Vougioukas, Y. Panagakis, S. Petridis, J. Kossaifi, and M. Pantic, "Speech-driven facial animation using polynomial fusion of features," in *Proc. IEEE-ICASSP*, 2020, pp. 3487–3491.

[261] S. Sinha, S. Biswas, and B. Bhowmick, "Identity-preserving realistic talking face generation," in *Proc. IEEE-IJCNN*, 2020, pp. 1–10.

[262] W. Wang, Y. Wang, J. Sun, Q. Liu, J. Liang, and T. Li, "Speech driven talking head generation via attentional landmarks based representation." in *Proc. Interspeech*, 2020, pp. 1326–1330.

[263] S. E. Eskimez, R. K. Maddox, C. Xu, and Z. Duan, "End-to-end generation of talking faces from noisy speech," in *Proc. IEEE-ICASSP*, 2020, pp. 1948–1952.

[264] R. Yi, Z. Ye, J. Zhang, H. Bao, and Y.-J. Liu, "Audio-driven talking face video generation with learning-based personalized head pose," *arXiv preprint arXiv:2002.10137*, 2020.

[265] L. Chen, G. Cui, C. Liu, Z. Li, Z. Kou, Y. Xu, and C. Xu, "Talking-head generation with rhythmic head motion," in *Proc. ECCV*, 2020.

[266] G. Mittal and B. Wang, "Animating face using disentangled audio representations," in *Proc. IEEE/CVF-WACV*, 2020, pp. 3290–3298.

[267] H. Zhu, H. Huang, Y. Li, A. Zheng, and R. He, "Arbitrary talking face generation via attentional audio-visual coherence learning," in *Proc. IJCAI*, 2021, pp. 2362–2368.

[268] C. Zhang, Y. Zhao, Y. Huang, M. Zeng, S. Ni, M. Budagavi, and X. Guo, "Facial: Synthesizing dynamic talking face with implicit attribute learning," in *Proc. IEEE/CVF-ICCV*, 2021, pp. 3867–3876.

[269] S. Si, J. Wang, X. Qu, N. Cheng, W. Wei, X. Zhu, and J. Xiao, "Speech2video: Cross-modal distillation for speech to video generation," *arXiv preprint arXiv:2107.04806*, 2021.

[270] S. Chen, Z. Liu, J. Liu, Z. Yan, and L. Wang, "Talking head generation with audio and speech related facial action units," *arXiv preprint arXiv:2110.09951*, 2021.

[271] 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 *Proc. IEEE/CVF-CVPR*, 2021, pp. 4176–4186.

[272] B. Liang, Y. Pan, Z. Guo, H. Zhou, Z. Hong, X. Han, J. Han, J. Liu, E. Ding, and J. Wang, "Expressive talking head generation with granular audio-visual control," in *Proc. IEEE/CVF-CVPR*, 2022, pp. 3387–3396.

[273] S. Wang, L. Li, Y. Ding, and X. Yu, "One-shot talking face generation from single-speaker audio-visual correlation learning," in *Proc. Conf AAAI Artif. Intell.*, vol. 36, no. 3, 2022, pp. 2531–2539.

[274] X. Ji, H. Zhou, K. Wang, Q. Wu, W. Wu, F. Xu, and X. Cao, "Eamm: One-shot emotional talking face via audio-based emotion-aware motion model," in *Proc. ACM SIGGRAPH*, 2022, pp. 1–10.

[275] S. Gururani, A. Mallya, T.-C. Wang, R. Valle, and M.-Y. Liu, "Spacex: Speech-driven portrait animation with controllable expression," *arXiv preprint arXiv:2211.09809*, 2022.

[276] R. Wu, Y. Yu, F. Zhan, J. Zhang, X. Zhang, and S. Lu, "Audio-driven talking face generation with diverse yet realistic facial animations," *arXiv preprint arXiv:2304.08945*, 2023.

[277] J. Liu, X. Wang, X. Fu, Y. Chai, C. Yu, J. Dai, and J. Han, "Opt: One-shot pose-controllable talking head generation," in *Proc. IEEE-ICASSP*, 2023, pp. 1–5.

[278] D. Wang, Y. Deng, Z. Yin, H.-Y. Shum, and B. Wang, "Progressive disentangled representation learning for fine-grained controllable talking head synthesis," in *Proc. IEEE/CVF-CVPR*, 2023, pp. 17 979–17 989.

[279] L. Zhang, Q. Chen, and Z. Liu, "Talking head generation for media interaction system with feature disentanglement," in *Proc. IEEE-ICPADS*. IEEE, 2023, pp. 403–410.

[280] O. Wiles, A. Koepke, and A. Zisserman, "X2face: A network for controlling face generation using images, audio, and pose codes," in *Proc. ECCV*, 2018, pp. 670–686.

[281] A. Jamaludin, J. S. Chung, and A. Zisserman, "You said that?: Synthesising talking faces from audio," *International Journal of Computer Vision*, vol. 127, pp. 1767–1779, 2019.

[282] X. Wen, M. Wang, C. Richardt, Z.-Y. Chen, and S.-M. Hu, "Photorealistic audio-driven video portraits," *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 12, pp. 3457–3466, 2020.

[283] A. Lahiri, V. Kwatra, C. Frueh, J. Lewis, and C. Bregler, "Lipsync3d: Data-efficient learning of personalized 3d talking faces from video using pose and lighting normalization," in *Proc. IEEE/CVF-CVPR*, 2021, pp. 2755–2764.

[284] Y. Lu, J. Chai, and X. Cao, "Live speech portraits: real-time photorealistic talking-head animation," *ACM Transactions on Graphics (TOG)*, vol. 40, no. 6, pp. 1–17, 2021.

[285] D. Bigioi, H. Jordan, R. Jain, R. McDonnell, and P. Corcoran, "Pose-aware speech driven facial landmark animation pipeline for automated dubbing," *IEEE Access*, vol. 10, pp. 133 357–133 369, 2022.

[286] W. Zhang, X. Cun, X. Wang, Y. Zhang, X. Shen, Y. Guo, Y. Shan, and F. Wang, "Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation," in *Proc. IEEE/CVF-CVPR*, 2023, pp. 8652–8661.

[287] S. Yao, R. Zhong, Y. Yan, G. Zhai, and X. Yang, "Dfa-nerf: personalized talking head generation via disentangled face attributes neural rendering," *arXiv preprint arXiv:2201.00791*, 2022.

[288] S. Shen, W. Li, Z. Zhu, Y. Duan, J. Zhou, and J. Lu, "Learning dynamic facial radiance fields for few-shot talking head synthesis," in *Proc. ECCV*, 2022, pp. 666–682.

[289] X. Liu, Y. Xu, Q. Wu, H. Zhou, W. Wu, and B. Zhou, "Semantic-aware implicit neural audio-driven video portrait generation," in *Proc. ECCV*, 2022, pp. 106–125.

[290] Z. Yu, Z. Yin, D. Zhou, D. Wang, F. Wong, and B. Wang, "Talking head generation with probabilistic audio-to-visual diffusion priors," *arXiv preprint arXiv:2212.04248*, 2022.

[291] Y. Zhua, C. Zhanga, Q. Liub, and X. Zhoub, "Audio-driven talking head video generation with diffusion model," in *Proc. IEEE-ICASSP*. IEEE, 2023, pp. 1–5.

[292] C. Xu, S. Zhu, J. Zhu, T. Huang, J. Zhang, Y. Tai, and Y. Liu, "Multimodal-driven talking face generation via a unified diffusion-based generator," *CoRR*, 2023.

[293] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: delving deep into convolutional nets," in *Proc. BMVC*. British Machine Vision Association, 2014.

[294] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for

unconstrained photo collections," in *Proc. IEEE/CVF-CVPR*, 2021, pp. 7210–7219.

[295] R. Huang, Z. Zhao, H. Liu, J. Liu, C. Cui, and Y. Ren, "Prodiff: Progressive fast diffusion model for high-quality text-to-speech," in *Proc. ACM MM*, 2022, pp. 2595–2605.

[296] C. Saharia, W. Chan, H. Chang, C. Lee, J. Ho, T. Salimans, D. Fleet, and M. Norouzi, "Palette: Image-to-image diffusion models," in *Proc. ACM SIGGRAPH*, 2022, pp. 1–10.

[297] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[298] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proc. IEEE-CVPR*, 2018, pp. 586–595.

[299] S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz, "Mocogan: Decomposing motion and content for video generation," in *Proc. IEEE-CVPR*, 2018, pp. 1526–1535.

[300] J. Chung and A. Zisserman, "Out of time: automated lip sync in the wild," in *Proc. ACCV*, 2017.

[301] N. D. Narvekar and L. J. Karam, "A no-reference perceptual image sharpness metric based on a cumulative probability of blur detection," in *Proc. IEEE-QoMEX*, 2009, pp. 87–91.

[302] K. De and V. Masilamani, "Image sharpness measure for blurred images in frequency domain," *Procedia Engineering*, vol. 64, pp. 149–158, 2013.

[303] D. Zeng, S. Zhao, J. Zhang, H. Liu, and K. Li, "Expression-tailored talking face generation with adaptive cross-modal weighting," *Neurocomputing*, vol. 511, pp. 117–130, 2022.

[304] K. Vougioukas, S. Petridis, and M. Pantic, "Realistic speech-driven facial animation with GANs," *International Journal of Computer Vision*, vol. 128, pp. 1398–1413, 2020.

[305] J. Wang, J. Liu, L. Zhao, S. Wang, R. Yu, and L. Liu, "Acoustic-to-articulatory inversion based on speech decomposition and auxiliary feature," in *Proc. IEEE-ICASSP*, 2022, p. 4808–4812.

[306] M. Mori, K. MacDorman, and N. Kageki, "The uncanny valley [from the field]," *IEEE Robotics and Automation Magazine*, vol. 19, pp. 98–100, 06 2012.

[307] C. Sheng, M. Pietikäinen, Q. Tian, and L. Liu, "Cross-modal self-supervised learning for lip reading: When contrastive learning meets adversarial training," in *Proc. ACM MM*, 2021, pp. 2456–2464.

[308] Y. Mroueh, E. Marcheret, and V. Goel, "Deep multimodal learning for audio-visual speech recognition," in *Proc. IEEE-ICASSP*, 2015, pp. 2130–2134.

[309] B. Shillingford, Y. Assael, M. W. Hoffman, T. Paine, C. Hughes, U. Prabhu, H. Liao, H. Sak, K. Rao, L. Bennett *et al.*, "Large-scale visual speech recognition," *Proc. Interspeech*, pp. 4135–4139, 2019.

[310] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *Proc. ICML*, 2021, pp. 8748–8763.

[311] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *Proc. ICML*, 2021, pp. 8821–8831.

[312] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs trained by a two time-scale update rule converge to a local nash equilibrium," *Proc. Advances in neural information processing systems (NIPS)*, vol. 30, 2017.

[313] E. Mu, K. M. Lewis, A. V. Dalca, and J. Guttag, "Generating image-specific text improves fine-grained image classification," *arXiv preprint arXiv:2307.11315*, 2023.