

SignBERT+: Hand-model-aware Self-supervised Pre-training for Sign Language Understanding

Hezhen Hu, Weichao Zhao*, Wengang Zhou, *Senior Member, IEEE*, and Houqiang Li, *Fellow, IEEE*

Abstract—Hand gesture serves as a crucial role during the expression of sign language. Current deep learning based methods for sign language understanding (SLU) are prone to over-fitting due to insufficient sign data resource and suffer limited interpretability. In this paper, we propose the *first* self-supervised pre-trainable SignBERT+ framework with model-aware hand prior incorporated. In our framework, the hand pose is regarded as a visual token, which is derived from an off-the-shelf detector. Each visual token is embedded with gesture state and spatial-temporal position encoding. To take full advantage of current sign data resource, we first perform self-supervised learning to model its statistics. To this end, we design multi-level masked modeling strategies (joint, frame and clip) to mimic common failure detection cases. Jointly with these masked modeling strategies, we incorporate model-aware hand prior to better capture hierarchical context over the sequence. After the pre-training, we carefully design simple yet effective prediction heads for downstream tasks. To validate the effectiveness of our framework, we perform extensive experiments on three main SLU tasks, involving isolated and continuous sign language recognition (SLR), and sign language translation (SLT). Experimental results demonstrate the effectiveness of our method, achieving new state-of-the-art performance with a notable gain.

Index Terms—Self-supervised pre-training, masked modeling strategies, model-aware hand prior, sign language understanding

1 INTRODUCTION

SIGN language (SL) serves as a primary communication tool for the deaf community. It is a visual language with its unique grammar and lexicon, which is non-trivial for the hearing people to master. To facilitate barrier-free communication between hearing and deaf people, automatic visual sign language understanding, as a topic with broad social influence, has been widely studied. Visual sign language understanding contains three main tasks, including isolated and continuous sign language recognition (SLR) and sign language translation (SLT). Isolated SLR focuses on word level recognition, which is essentially a fine-grained classification task. Differently, continuous SLR targets at recognizing the sign gloss sequence with its corresponding occurring order, which needs to learn the sequence correspondence across the visual and textual domains. SLT intends to further generate spoken language translations, which emphasizes natural linguistic expression. These three tasks are all important for sign language understanding and bring challenges from different perspectives.

Hand gesture plays a dominant role during the meaning expression of sign language. Intrinsically, hand occupies a relatively small spatial size, exhibiting uniform appearance and self-occlusion among joints. During SL expression, it usually occurs over complex backgrounds and presents

fast motion. These characteristics lead to difficulty in hand representation learning. Current deep-learning-based methods adaptively learn hand feature representations from the cropped RGB sequence. Meanwhile, considering the highly-articulated characteristic of hand, some methods propose to utilize pose as the hand representation. Pose is a compact and semantic representation, which is robust to appearance change and brings potential computation efficiency. However, current pose-based methods utilize the poses extracted from off-the-shelf pose detectors, which usually suffer failure detection due to the motion blur and complex backgrounds, *etc.* Therefore, their performances usually lag largely behind the RGB-based counterparts. Besides, the aforementioned methods all follow a data-driven paradigm and suffer over-fitting due to limited sign data resource and insufficient interpretability.

Meanwhile, the effectiveness of pre-training has been validated in computer vision (CV) and natural language processing (NLP) tasks. Recent advances in NLP are largely derived from self-supervised pre-training techniques on large text corpus [1], [2], [3]. Among them, BERT [2] is one of the most popular methods for its simplicity and effectiveness. Its success is largely attributed to the strong Transformer backbone [4] and well-designed pre-training strategies to model the context inherent in the text corpus. By adding a simple head (an MLP) on top for fine-tuning, it achieves notable performance gains in many downstream tasks, especially those with limited data resource. Notably, natural language is represented by a 1D sequence of text words which are characterised with well-defined semantic meaning. However, sign video is a kind of 3D data with complex spatial-temporal context, and it is non-trivial to analogically define visual word or unit with clarified semantics. Therefore, it remains a hard challenge to leverage the success of BERT to video-based sign language understand-

- This work was supported by the National Natural Science Foundation of China under Contract U20A20183 and 62021001. It was also supported by GPU cluster built by MCC Lab of Information Science and Technology Institution, USTC, and the Supercomputing Center of the USTC.
- Hezhen Hu, Weichao Zhao, Wengang Zhou and Houqiang Li are with the Department of Electronic Engineering and Information Science of Electrical and Computer Engineering, University of Science and Technology of China, Hefei, Anhui, 230026.
E-mail: {alexhu, saruka}@mail.ustc.edu.cn, {zhtwg, lihq}@ustc.edu.cn
*Equal contribution with the first author.
Corresponding authors: Wengang Zhou and Houqiang Li.

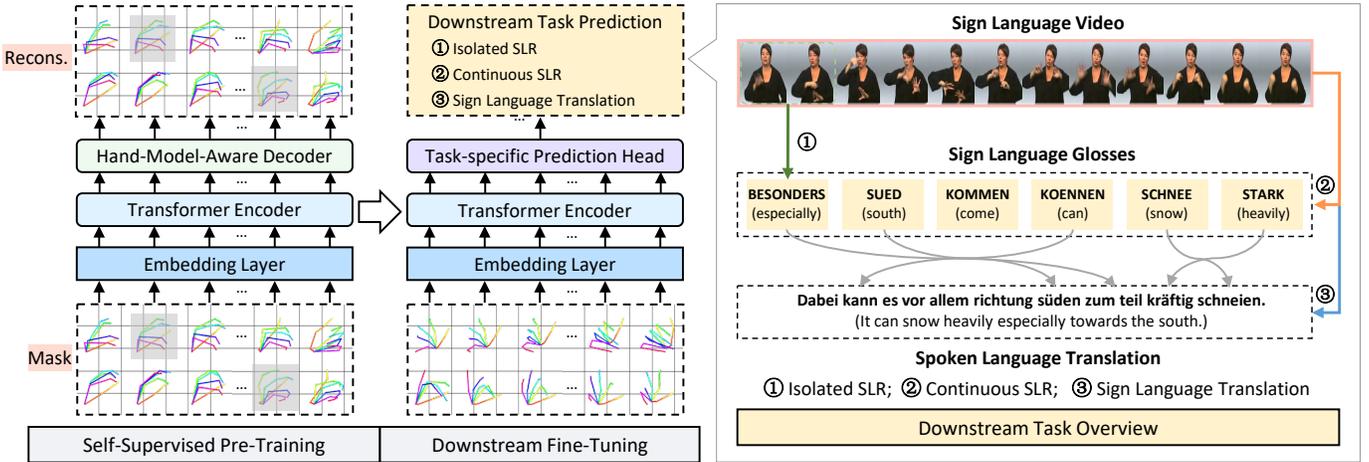


Fig. 1. The overview of our method and sign language understanding tasks (isolated SLR, continuous SLR and SLT).

ing.

To tackle the above-mentioned issues, we develop a self-supervised pre-trainable framework with model-aware hand prior incorporated, namely SignBERT+, as shown in Figure 1. Considering the dominance of hand during SL expression, we utilize the compact and expressive hand pose as a visual token in a frame-by-frame manner. Then, we carefully depict it with gesture state and spatial-temporal global position information. SignBERT+ first performs self-supervised pre-training on a large volume of hand pose data, which is derived from an off-the-shelf detector on sign videos. Inspired by the success of BERT [2], we pre-train the encoder-decoder backbone via reconstructing the masked visual tokens from the corrupted input sequence, which enforces the framework to capture the hierarchical context in the sign language domain. Considering the noisy characteristic of detected hand pose data, we carefully design multi-scale masking strategies, including joint, frame and clip levels. Meanwhile, to better mine the context in the sign video domain, we further incorporate hand prior in a model-aware method. After pre-training, we carefully design simple yet effective task-specific prediction heads, which are jointly fine-tuned with the pre-trained SignBERT+ encoder to adapt to downstream tasks.

In summary, our contributions are three-fold as follows,

- To our best knowledge, we propose the *first* model-aware pre-trainable framework, namely SignBERT+. It performs self-supervised pre-training on a large volume of sign pose data, followed by fine-tuning to achieve better performance on multiple downstream tasks.
- To better model the hierarchical context underneath the sign data during pre-training, we design multiple masked modeling strategies ranging from joint to clip level, in coordination with incorporated model-aware hand prior and spatial-temporal position encoding. For diverse downstream tasks, we design simple yet effective task-specific prediction heads on top of the pre-trained SignBERT+ encoder.
- We perform extensive experiments to validate the feasibility and effectiveness of our framework. Experimental results demonstrate that our method

achieves new state-of-the-art performance on video-based sign language understanding tasks, including isolated SLR, continuous SLR and SLT.

This work is an extension of the conference paper [5] with improvement in a number of aspects. 1) Considering the characteristics of sign language, we further introduce spatial-temporal global position encoding into embedding, along with the masked clip modeling for modeling temporal dynamics. Those new techniques further bring a notable performance gain. 2) We extend the original framework to two more downstream tasks in video-based sign language understanding, *i.e.*, continuous SLR and SLT. To this end, we design simple yet effective task-specific prediction heads. Besides, we also provide efficient fusion strategies with the RGB modality. Our newly designed framework achieves state-of-the-art performance on all the downstream tasks. 3) We present more comprehensive discussion on related works and make deep analysis on different components of our method to highlight the important ingredients. Besides, we add discussions on future works and broader impact.

2 RELATED WORK

In this section, we first give a literature review for video-based sign language understanding. Then we present an overview of pre-training strategies. Finally, we introduce related hand modeling techniques.

2.1 Video-based Sign Language Understanding

Video-based sign language understanding has made remarkable progress [6], [7], [8], [9]. Generally, it contains three main tasks, including isolated SLR, continuous SLR and SLT. These tasks emphasize different aspects, bringing their specific challenges to resolve.

Isolated sign language recognition. Isolated SLR aims to recognize at the word level, which is essentially a fine-grained classification problem. This task poses a challenge on learning discriminative visual representation [7], [10], [11], [12], [13]. Early works utilize hand-crafted features, *e.g.* HOG [14] and SIFT [10], to represent hand shape, orientation and motion. Recently, researchers have resorted to deep learning techniques, which adaptively extract features

from the full video sequence. Based on the input modality, these works can be divided into RGB-based and pose-based methods. RGB-based methods usually adopt Convolutional Neural Networks (CNNs) as the backbone. For instance, Koller *et al.* [15] utilize 2D-CNNs with LSTM to sequentially model the spatial and temporal representations. Some other works utilize 3D-CNNs for modeling spatial-temporal dependency [7], [13], [16], [17], [18].

For the pose-based counterpart, there exist different backbones for feature extraction, including CNNs [13], [19] and RNNs [20], [21], [22], *etc.* Recently, considering its well-structured nature, more and more works have utilized graph convolutional networks (GCNs), which exhibit both efficiency and effectiveness [20], [22], [23]. As a representative work, ST-GCN [24] organizes the pose sequence as a pre-defined graph and adopts GCNs to perform recognition. Besides, Tunga *et al.* further combines Transformer without pre-training for isolated SLR [23].

Continuous sign language recognition. It aims to map the sign video to the gloss sequence in the same presenting order. In this task, the transitions between sign glosses may come with temporal variants, and the sign video usually lacks the frame-level gloss annotation. Therefore, it raises a new challenge on the sequence correspondence learning between the visual sign representation to the sign glosses. To this end, Koller *et al.* [25], [26] exploit the integration of 2D-CNNs and Hidden Markov Models (HMMs) for modeling transitions. Connectionist Temporal Classification (CTC) is a differentiable cost function, which is able to deal with two unsegmented sequences without precise alignment. It usually works with Recurrent Neural Networks (RNNs), *e.g.* BLSTM [27] and GRU [28], and Transformer [4] for sequential learning. CTC-based methods make end-to-end optimization possible and become the mainstream for its competitive performance [8], [9], [29], [30], [31]. However, these methods are prone to over-fitting due to limited data resources. To tackle this issue, DNF [29] utilizes the iterative optimization strategy for better feature representation. Zhou *et al.* [32] boosts the visual encoder with the partially masked videos under the supervised classification task. CMA [30] proposes cross modality augmentation, which leverages the pseudo video-text pairs to boost recognition performance. VAC [8] proposes visual alignment constraint to enhance the feature extractor.

Sign language translation. This task intends to generate the spoken language translations. It is mainly different from continuous SLR in the aspect of sequential learning, due to different grammar and word order between sign language and spoken language [33], [34]. NSLT [33] first explores this task with an attention-based encoder-decoder and proposes RWTH-PhoenixT dataset. This dataset provides both sign gloss and translation annotation, and becomes the most popular benchmark. Camgoz *et al.* [9] leverage the strong modeling capability of Transformer into sequential learning. TSPNet [34] explores the temporal semantic structures for more discriminative features. STMC [35] fuses information from multi-cue streams to boost performance. SignBT [36] utilizes external text corpus for performance improvement. In this work, we aim to leverage a large volume of sign data via pre-training to benefit three main sign language understanding tasks.

2.2 Pre-Training Strategy

Pre-training, as a common strategy in CV and NLP, aims to learn generic representation from massive labeled or unlabeled data, which benefits downstream tasks with marginal fine-tuning cost. For fully supervised pre-training, it is common for CV tasks to first pre-train CNNs under labeled classification benchmarks, *e.g.* ImageNet [37] and Kinetics [38], *etc.* However, given the labeling cost, more and more works turn to self-supervised learning from a large volume of unlabeled data, which is readily available from the Web [39], [40]. Self-supervised learning aims to model the joint probability distribution inherent in data, which is beneficial to address the following discriminative learning task.

Pioneering works subtly design pretext tasks to perform self-supervised pre-training [41], [42], [43], [44], [45], [46], [47]. These tasks include predicting colorization [41], rotation [43], transformation [46] and frame / clip orders [47], [48], *etc.* Recently, some works focus on contrastive learning for pre-training [49], [50], [51]. Typically, it aims to pull the representation of similar instances closer, while pushing away negative instances. To obtain informative negative instances for better optimization, some works utilize the techniques like memory banks [49] and large batch size [52]. There also exist works further eliminating the requirement of negative samples [51], [53].

Another interesting strand is the generative self-supervised pre-training, which usually involves training the encoder via the reconstruction task. In NLP, one milestone of pre-training is BERT [2]. BERT is built on the strong Transformer backbone with masked language modeling (MLM) as one of its pre-training tasks. MLM attempts to predict the masked words by leveraging the context cues from the remaining tokens. Similar to BERT, some works also adopt MLM during pre-training, such as GPT [1], XLNet [3] and RoBERTa [54], *etc.* These pre-training methods generalize well and bring notable performance gains on the downstream tasks. Motivated by the success in NLP, some works attempt to leverage the idea of BERT into CV tasks [55], [56], [57], [58], [59], [60], which mainly focus on the RGB modality. BEiT [59] utilizes the discrete tokenized image patches as pseudo labels and performs masked modeling similar to BERT. MAE [60] directly works in the continuous space, *i.e.*, masking and reconstructing the pixel values. However, BEiT and MAE only focus on image-based tasks. Actually, it is non-trivial to leverage BERT's success to video-based sign language understanding tasks, which involves special design of the pretext task and framework architecture.

Pre-training in sign language. Albanie *et al.* [13] propose to perform supervised pre-training on a large-scale annotated dataset. Li *et al.* [7] boost isolated SLR via a domain-invariant feature descriptor, which leverages the knowledge from external subtitled news sign video. To our best knowledge, there exists no self-supervised pre-training works in the sign language domain.

2.3 Hand Modeling Techniques

Hand modeling aims to depict the hand with more expressiveness. Current modeling techniques include sum-of-Gaussians [61], shape primitives [62], [63] and sphere-

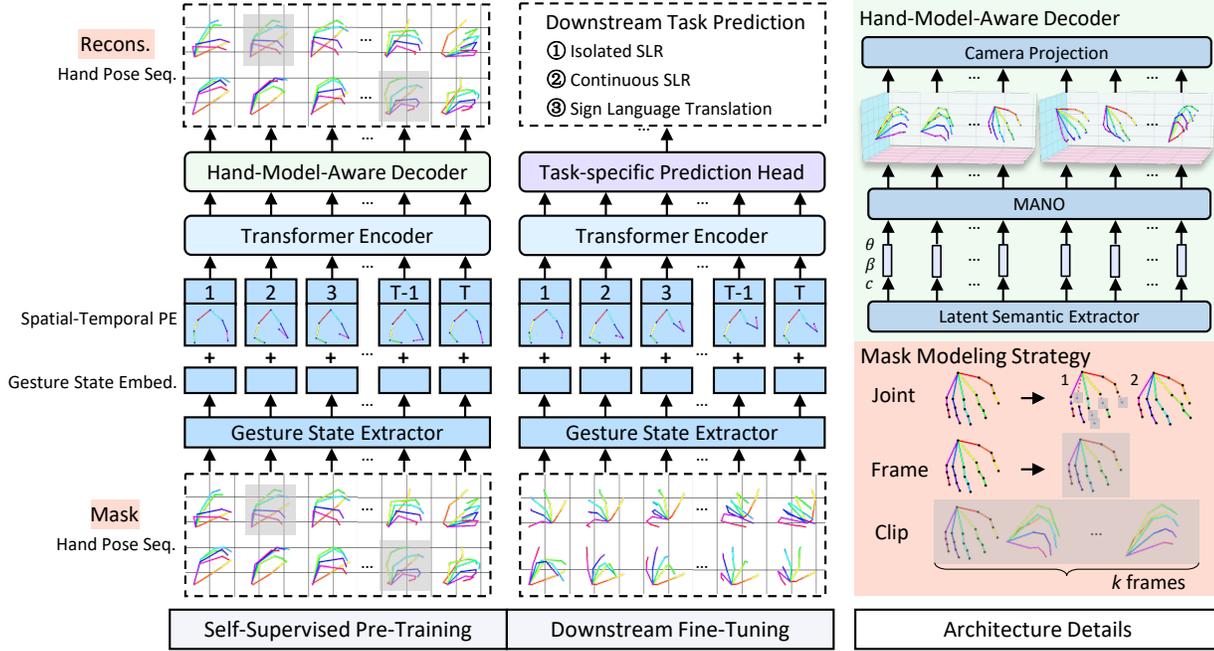


Fig. 2. Illustration of our SignBERT+ framework details, which contains self-supervised pre-training and fine-tuning for the downstream tasks. We organize the pre-extracted 2D poses of both hands as the visual token sequence. For each token, it is embedded with gesture state and spatial-temporal position encoding. During self-supervised pre-training, multi-level masked modeling strategies work with incorporated model-aware hand prior, in order to better capture the hierarchical context in the sign domain. Given the downstream-task diversity, we design the task-specific prediction head and fine-tune it with the pre-trained SignBERT+ encoder.

meshes [64], *etc.* To better reconstruct the hand shape, Iason *et al.* [65] define scaling terms on bone lengths. Notably, some works [66], [67], [68] attempt to learn hand shape variation with Linear Blend Skinning (LBS) [69]. Among them, MANO [68] becomes the most popular one for its wide applications [70], [71], [72], [73]. As a statistical model, it learns from a large variety of high-quality hand scans and represents the geometric changes in the low-dimensional pose and shape space. With this capability, we adopt it as a constraint in the decoder to incorporate prior.

3 OUR APPROACH

As shown in Figure 2, our framework contains two stages, *i.e.*, self-supervised pre-training and downstream task fine-tuning. During sign expression, both hands are involved to act as a dominant role. Therefore, we focus on them to build the visual token in a frame-wise manner. For each visual token, we embed the gesture state and global spatial-temporal position information. During pre-training, the whole framework works in a self-supervised paradigm by reconstructing the masked visual tokens from the corrupted input sequence. Jointly with the multi-level masking strategies, we incorporate hand prior to better capture hierarchical context in the sign domain. Then, we fine-tune the pre-trained SignBERT+ encoder (embedding layer and Transformer encoder) with the designed prediction heads for downstream tasks.

In the following, we first introduce the framework architecture. After that, we elaborate our pre-training strategy. Finally, we discuss the fine-tuning schemes for downstream tasks.

3.1 Framework Architecture

Our framework contains four main components, *i.e.*, input embedding layer, Transformer encoder, hand-model-aware decoder and prediction head.

3.1.1 Input Embedding Layer

Given the dominant role of hand during sign language, we carefully design an embedding layer to capture cues from both hands. It extracts gesture state and spatial-temporal position information from the hand pose sequence in a frame-wise manner, which are elaborated in the follow paragraphs.

Gesture state embedding. Given its well-structured characteristics, we first organize the input 2D hand pose \tilde{J}_t at frame t as an undirected spatial graph. This graph is constructed with the node V and edge E set. The node set includes all hand joints, *i.e.*, 21 joints per hand, while the edge set contains their physical and symmetrical connection. Then, the hand pose sequence is processed by a spectral-based GCN [24], [74], which hierarchically performs graph convolution and graph pooling for gesture state embedding. The graph convolution is formulated as follows,

$$f_{out} = \sum_i D_i^{-\frac{1}{2}} A_i D_i^{\frac{1}{2}} f_{in} W_i, \quad (1)$$

where f_{in} and f_{out} are the corresponding input and output features, respectively. i indicates the type of neighbors for each node. W_i denotes the convolution weight and A_i is the dismantled matrix indicating the edge connection. For graph pooling, we first cluster the original 21 joint nodes of each hand into 6 subsets corresponding to 5 fingers and 1 palm. Then we perform max-pooling on the nodes in each

subset, leading to 6 nodes. Finally, these nodes are again max-pooled into one and both hands are involved in the frame-level gesture state embedding $\mathbf{f}_{p,t}$.

Spatial-temporal position encoding. Besides the gesture state, hand spatial trajectory and temporal information also matter in video-based sign language understanding. We depict the hand global position in the normalized 2D space by introducing the arm joints of both sides. These joints are also processed by GCN [24], [74] to extract the frame-level spatial embedding $\mathbf{f}_{s,t}$. Since the Transformer layers process the sequence in an order-agnostic way, we add temporal information into the input embedding $\mathbf{f}_{e,t}$, which is implemented by the position encoding technique in [4].

3.1.2 Transformer Encoder

The embedded input sequence is fed into the Transformer encoder [4] for the latent semantic representation. Its basic layer mainly contains two components, *i.e.*, a multi-head self-attention module and a feed-forward network. For each layer, its output retains the same size with the input. The whole encoder is formulated as follows,

$$\begin{aligned} \mathbf{F}_0 &= \{\mathbf{f}_{p,t} + \mathbf{f}_{s,t} + \mathbf{f}_{e,t}\}_{t=1}^T, \\ \tilde{\mathbf{F}}_i &= L(M(\mathbf{F}_{i-1}) + \mathbf{F}_{i-1}), \\ \mathbf{F}_i &= L(C(\tilde{\mathbf{F}}_i) + \tilde{\mathbf{F}}_i), \end{aligned} \quad (2)$$

where i denotes the i -th layer of the Transformer encoder. The whole encoder contains totally N layers. $M(\cdot)$, $C(\cdot)$ and $L(\cdot)$ represent the multi-head self-attention, feed-forward network and layer normalization, respectively. \mathbf{F}_i denotes the output feature from the i -th layer.

3.1.3 Hand-model-aware Decoder

To achieve the reconstruction target during pre-training, the decoder transforms the latent feature back to the pose sequence. The decoder works in a model-aware method to incorporate prior, which aims to guide the encoder better capturing generic representations in the sign language domain. Specifically, the latent feature is first processed by a fully-connected layer, which extracts the low-dimensional semantic embeddings depicting the hand status, *i.e.*, hand pose θ and shape β , and the camera parameter c aligning the image plane, which is formulated as follows,

$$\mathbf{F}_{la} = \{\theta, \beta, c_r, c_o, c_s\}_{t=1}^T = D(\mathbf{F}_N), \quad (3)$$

where $D(\cdot)$ denotes the fully-connected layer. θ and $\beta \in \mathbb{R}^{10}$ are the hand pose and shape embeddings for the following MANO model, respectively. $c_r \in \mathbb{R}^{3 \times 3}$, $c_o \in \mathbb{R}^2$, and $c_s \in \mathbb{R}$ are parameters of the weak-perspective camera, representing the rotation, translation and scale, respectively.

Then the MANO model [68] incorporates hand prior and decodes the estimated hand embedding. Specifically, the decoding process is fully-differentiable, which transforms the hand embedding (θ and β) to the dense triangular hand mesh $\mathbf{M} \in \mathbb{R}^{N_v \times 3}$ ($N_v = 778$ vertices and $N_f = 1538$ faces) as follows,

$$\mathbf{M}(\beta, \theta) = W(\mathbf{T}(\beta, \theta), J(\beta), \theta, \mathbf{W}), \quad (4)$$

$$\mathbf{T}(\beta, \theta) = \bar{\mathbf{T}} + B_S(\beta) + B_P(\theta), \quad (5)$$

where $B_S(\cdot)$ and $B_P(\cdot)$ represent shape and pose blend functions, respectively. \mathbf{W} is a set of blend weights. The

hand template $\bar{\mathbf{T}}$ is first posed and skinned based on the pose and shape corrective blend shapes, *i.e.*, $B_P(\theta)$ and $B_S(\beta)$. Then the mesh \mathbf{M} is generated by rotating each part around joints $J(\beta)$ based on the linear skinning function $W(\cdot)$ [75]. Besides, the sparse joint representation \tilde{J}_{3D} is also derived from the mesh. To keep consistent with the input pose format, we further add 5 extra fingertip joints by selecting vertices with the index of 333, 443, 555, 678 and 734. Finally, \tilde{J}_{3D} is mapped back to the same 2D plane as the input pose based on the estimated camera parameter as follows,

$$\tilde{J}_{2D} = c_s \prod (\mathbf{c}_r \tilde{J}_{3D}) + \mathbf{c}_o, \quad (6)$$

where $\prod(\cdot)$ denotes the orthographic projection.

3.1.4 Prediction Head

Given the large diversities among downstream tasks, we design simple yet effective prediction heads for each task in Figure 3. In Section 3.3, we will introduce them in detail along with the task-specific fine-tuning settings.

3.2 Pre-Training SignBERT+

In this section, we elaborate our pre-training paradigm. Pre-training is performed via reconstructing the masked visual tokens from the corrupted input sequence, which aims to exploit hierarchical context on a large volume of sign pose data. Different from the original BERT working on discrete word space, we attempt to pre-train on continuous pose space. Therefore, it raises new issues to resolve, including the design of the masking strategies and objective functions.

3.2.1 Masking Strategy

Considering the noise of the detected input hand pose, the masking strategy needs to be carefully redesigned. Given the hand pose sequence, we first randomly choose a portion R of all tokens. For the chosen token, one of the following operations is applied with the equal probability, *i.e.*, masked joint modeling, masked frame modeling, masked clip modeling and identity modeling.

Masked joint modeling. This strategy mimics the failure cases of pose detectors on some joints. For a chosen token, we randomly choose m joints. Two operations are performed on these chosen joints with the equal probability, *i.e.*, zero masking (masking the coordinates of joints with zeros) or random spatial disturbance. This modeling aims to guide the framework to infer the gesture state from the remaining joints, thus capturing the context at the joint level.

Masked frame modeling. It aims to deal with the failure case on the whole frame pose, which is often caused by complex backgrounds. The chosen token is directly zero masked. This strategy enforces the framework to reconstruct the token by leveraging the observation from the remaining frames and the other hand. In this way, the temporal context for each hand and mutual context between hands are captured.

Masked clip modeling. Motion blur, as a factor not to be overlooked, usually causes pose detection failure on a video clip. To deal with this situation, masked clip modeling is designed. We randomly choose k temporally continuous tokens, where k ranges from 2 to K . The chosen k tokens

are all zero-masked. In order to reconstruct them, the framework needs to capture the temporal dynamics by leveraging the motion pattern of existing frames.

Identity modeling. Similar to BERT [2], identity modeling directly feeds the unchanged tokens into the framework. It is indispensable for the framework to learn identity semantic encoding on those unmasked tokens.

3.2.2 Objective Functions

During pre-training, its objective is to maximize the likelihood of the joint probability distribution to reconstruct the hand pose sequence. To achieve the reconstruction target, the classification objective in the original BERT is substantially changed into regression. To this end, we design the objective function as follows,

$$\mathcal{L} = \mathcal{L}_{rec} + \lambda \mathcal{L}_{reg}, \quad (7)$$

where \mathcal{L}_{rec} and \mathcal{L}_{reg} are reconstruction and regularization loss terms, respectively. λ denotes the weighting factor. We only include the corresponding output of the masked tokens during the loss calculation.

Reconstruction loss \mathcal{L}_{rec} . Since the utilized pose usually contains noise due to failure detection, we utilize the detection confidence score as the filter to eliminate these influences. The reconstruction loss is calculated as follows,

$$\mathcal{L}_{rec} = \sum_{t,j} \mathbb{1}(s(t,j) \geq \epsilon) s(t,j) \left\| \tilde{J}_{2D}(t,j) - J_{2D}(t,j) \right\|_1, \quad (8)$$

where $\mathbb{1}(\cdot)$ denotes the indicator function, and $s(t,j)$ denotes the detection confidence score of the J_{2D} with joint j at time t . The joints with the confidence lower than ϵ are not included in the loss calculation.

Regularization loss \mathcal{L}_{reg} . To ensure this decoder working properly, the regularization loss is added. It is implemented by constraining the magnitude and derivative of the MANO input, which is responsible for generating the plausible mesh and keeping the signer identity unchanged. This loss term is computed as follows,

$$\mathcal{L}_{reg} = \sum_t (\|\theta_t\|_2^2 + w_\beta \|\beta_t\|_2^2 + w_\delta \|\beta_t - \beta_{t-1}\|_2^2), \quad (9)$$

where w_β and w_δ denote the weighting factors.

3.3 Fine-Tuning SignBERT+

After pre-training SignBERT+ for generic visual representation in sign language, it is relatively simple to fine-tune it for various downstream tasks. During fine-tuning, the task-specific prediction head is added on top of the pre-trained SignBERT+ encoder, as illustrated in Figure 3. The framework is supervised under the task-specific loss.

Since only the hand pose modality is insufficient to convey the full meaning of sign language, we further provide the task-specific fusion strategy with the method based on the full RGB frames. For clarity, the baseline RGB method utilized for fusion will be marked in the experiment section. Besides, we denote our vanilla and fused framework as **Ours** and **Ours (+ R)**, respectively.

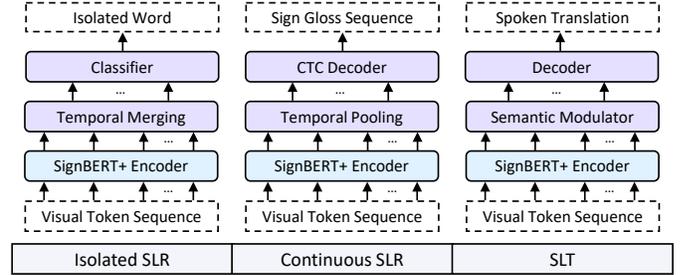


Fig. 3. Illustration of the settings on three downstream tasks, *i.e.*, isolated SLR, continuous SLR and SLT. The box in purple denotes our designed task-specific prediction head. It is fine-tuned with the pre-trained SignBERT+ encoder.

3.3.1 Isolated Sign Language Recognition

Isolated SLR is a fine-grained classification problem, which categorizes a sign video to the corresponding isolated word. For this task, the designed prediction head consists of a temporal merging module and a classifier. The former module utilizes a simple attention mechanism to highlight the discriminative cues in certain frames during the merging process as follows,

$$\mathbf{o} = \text{Softmax}(\mathbf{F}\mathbf{W} + b) \cdot \mathbf{F}, \quad (10)$$

where \mathbf{F} and \mathbf{o} denote the input feature sequence and merged feature, respectively. Then the merged feature \mathbf{o} is passed through a classifier (MLP and softmax layer) to output the probability matrix. Since isolated SLR is a classification problem, we utilize the cross-entropy loss to supervise the fine-tuning process. We use the simple late fusion strategy with the RGB method, which directly sums their prediction scores and chooses the class with the highest score as the final recognition result.

3.3.2 Continuous Sign Language Recognition

Continuous SLR aims to recognize the gloss sequence \mathbf{g} in the same presenting order as the sign actions in the input sign video \mathbf{V} with T frames. The prediction head for this task contains a temporal pooling module and a CTC decoder. The temporal pooling module aggregates frame-level visual features to the clip level, which outputs the one-quarter temporal length of the input. Then it is fed into the connectionist temporal classification (CTC) decoder to deal with the mapping between two unsegmented sequences without explicit alignment.

The objective of CTC is to maximize the posterior probability over all alignments from the source to the target. It extends the vocabulary with a blank label to cover the cases of transition and silence. Denote each alignment path of the input sequence as $\pi = \{\pi_t\}_{t=1}^T$ with T as temporal duration. Under the time independence assumption, its probability is computed as follows,

$$p(\pi|\mathbf{V}) = \prod_{t=1}^T p(\pi_t|\mathbf{V}). \quad (11)$$

Typically, there exists many-to-one mapping from multiple input sequences to one target, which is achieved by removing all blanks and repetition. In this way, we calculate the conditional probability of the target gloss sequence \mathbf{g} by

summing the probabilities of all possible mapping paths as follows,

$$p(\mathbf{g}|\mathbf{V}) = \sum_{\pi \in \mathcal{B}^{-1}(s)} p(\pi|\mathbf{V}), \quad (12)$$

where $\mathcal{B}(\cdot)$ denotes the many-to-one mapping function. $\mathcal{B}^{-1}(\cdot)$ is the inverse mapping of $\mathcal{B}(\cdot)$. During training, the objective is defined by the negative log probability of $p(\mathbf{g}|\mathbf{V})$ as follows,

$$\mathcal{L}_{\text{CSLR}} = -\ln p(\mathbf{g}|\mathbf{V}). \quad (13)$$

During inference, the CTC decoder obtains a series of sentences via beam search and chooses the one with the highest decoding probability as the final prediction.

For fusion, similar to [76], we first concatenate the encoded feature from the RGB baseline and our method. Then we utilize a BLSTM sequential model and a CTC decoder to map the merged feature to the gloss sequence.

3.3.3 Sign Language Translation

Given the input sign video \mathbf{V} with T frames, SLT aims to generate the spoken language translation $s = \{s_i\}_{i=1}^N$ with N words via maximizing the conditional probability $p(s|\mathbf{V})$. For this task, our designed prediction head contains a semantic modulator and a decoder as shown in Figure 3.

Considering the token length diversity between the source and target ($T \gg N$), the semantic modulator attempts to bridge this gap and generates suitable semantics $\mathbf{M} = \{\mathbf{m}_i\}_{i=1}^{T_1}$ for the decoder. Specifically, it first performs average temporal pooling to reduce the source visual token sequence from the length T to the length $T_1 = T/4$. This operation makes the visual representation more compact, but its output lacks temporal dependency modeling. To mitigate this issue, a Transformer encoder is adapted to further modulate the pooled visual sequence and generate suitable semantics.

After that, a decoder is adopted to perform mapping between sign language and spoken translation while considering their different grammar. The decoder contains two main components, *i.e.*, a word embedding layer and an autoregressive Transformer decoder. The word embedding layer embeds each word in the target sequence s , along with the added position encoding as follows,

$$\mathbf{w}_i = \text{WordEmbed}(s_i) + PE(i), \quad (14)$$

where s_i denotes the input word, $\text{WordEmbed}(\cdot)$ and $PE(\cdot)$ are the word embedding and position encoding functions, respectively.

The autoregressive property means the model leverages generated text as additional input when generating the next. The Transformer decoder architecture is also a stack of basic blocks. The basic block contains three components, *i.e.*, masked multi-head self-attention module, multi-head cross-attention module and feed-forward network. The mask adopted on self-attention ensures the information flow in the rightward direction to preserve the autoregressive property [4], [77]. This operation is necessary for the SLT inference, since the framework is not accessible to the output tokens which are decoded currently or in the future. Cross-attention module leverages the contextual cues from the

modulated visual semantics \mathbf{M} and predecessors words \mathbf{w} . The whole Transformer decoder is formulated as follows,

$$\begin{aligned} \mathbf{D}_0 &= \mathbf{w}, \\ \mathbf{Q}_i &= L(\widetilde{\mathbf{M}}(\mathbf{D}_{i-1}) + \mathbf{D}_{i-1}), \\ \widetilde{\mathbf{D}}_i &= L(MHA(\mathbf{Q}_i, \mathbf{M}^k, \mathbf{M}^v) + \mathbf{Q}_i), \\ \mathbf{D}_i &= L(C(\widetilde{\mathbf{D}}_i) + \widetilde{\mathbf{D}}_i), \end{aligned} \quad (15)$$

where i denotes the i -th layer of the Transformer decoder. The whole encoder contains totally N layers. $\widetilde{M}(\cdot)$, $MHA(\cdot)$, $C(\cdot)$ and $L(\cdot)$ represent the masked multi-head self-attention, multi-head cross-attention, feed-forward network and layer normalization, respectively. \mathbf{D}_i denotes the output feature from the i -th layer.

During decoding, the sentence is first prefixed with the word “[bos]” to indicate the beginning. Then each word in the target sequence s is embedded. The embedded sequence is then fed into the Transformer decoder $\text{TransD}(\cdot)$. This decoder additionally performs cross attention by leveraging the contextual cues from the modulated visual semantics \mathbf{m} and predecessors words \mathbf{w} . Finally, its output is fed into a fully-connected network and a softmax layer to generate the probability matrix of the output word.

In summary, the whole decoding process is formulated as follows,

$$\mathbf{h}_i = \text{TransD}(\mathbf{w}_{1:i-1}, \mathbf{m}_{1:T_1}), \quad (16)$$

$$\mathbf{o}_i = \text{Softmax}(\mathbf{W}\mathbf{h}_i + b). \quad (17)$$

The conditional probability $p(s|\mathbf{V})$ is calculated as follows,

$$p(s|\mathbf{V}) = \prod_{i=1}^N p(s_i | s_{1:i-1}, \mathbf{V}) = \prod_{i=1}^N \mathbf{o}_{i, s_i}. \quad (18)$$

Finally, the objective function is formulated as follows,

$$\mathcal{L}_{\text{SLT}} = -\ln p(s|\mathbf{V}), \quad (19)$$

which is equivalent to calculating the cross-entropy loss on each word. We adopt the S2T setting [33], which directly maps the sign embedding to spoken translation in an end-to-end manner. During inference, the framework predicts the word one-by-one based on the beam search [78].

For fusion, we leverage the latent visual features from both RGB and pose modalities, *i.e.*, \mathbf{M}_r and \mathbf{M}_p , to the same decoder. Specifically, we replace the original cross attention in the decoder with the cascaded one, which is formulated as follows,

$$\widetilde{\mathbf{D}}_i = MHA(MHA(\mathbf{Q}_i, \mathbf{M}_p^k, \mathbf{M}_p^v), \mathbf{M}_r^k, \mathbf{M}_r^v), \quad (20)$$

where $MHA(\cdot)$ denotes the multi-head cross-attention layer, \mathbf{Q}_i denotes the feature from the previous decoder layer, and $\widetilde{\mathbf{D}}_i$ is the output of the cascaded attention.

4 EXPERIMENT

In this section, we first introduce the experiment setup, *i.e.*, datasets, evaluation metrics and implementation details. Then we perform ablation studies on the framework effectiveness from multiple perspectives. Finally, we perform extensive experiments to make comparison with state-of-the-art methods on multiple downstream tasks.

4.1 Experiment Setup

4.1.1 Datasets

We first perform experiments on the dataset with hand pose annotations available to evaluate the framework feasibility. HANDS17 [79] is a video-level hand pose estimation dataset with 292,820 frames from 99 video sequences. For each video, the first 70% and remaining 30% frames are separated for training and testing, respectively.

We evaluate our proposed method on three main video-based sign language understanding tasks, *i.e.*, isolated SLR, continuous SLR and SLT. The corresponding datasets for each task are discussed in the follow.

For *isolated SLR*, we make evaluation on three datasets, *i.e.*, MSASL [17], WLASL [18], and SLR500 [16]. MSASL [17] is an American sign language (ASL) dataset containing a vocabulary of 1,000, with 25,512 samples. Besides, it also releases two subsets (MSASL100 and MSASL200) by choosing the Top- K most frequent signs. WLASL [18] is another ASL dataset with 2,000 signs and 21,083 samples. It also contains two subsets, *i.e.*, WLASL100 and WLASL300. MSASL and WLASL are both collected from the Web, which are recorded in unconstrained real-life conditions with unbalanced samples for each sign word. These factors bring new challenges on accurate recognition. SLR500 [16] is the largest CSL dataset, which contains 500 daily signs and 125,000 samples recording at the resolution of 1280×720 . These samples are split into 90,000 and 35,000 for training and testing, respectively.

For *continuous SLR*, the evaluation is conducted on two datasets, *i.e.*, RWTH-Phoenix [14] and RWTH-PhoenixT [33]. RWTH-Phoenix [14] is a popular German sign language dataset collected from the weather forecast broadcast. It contains 6,841 samples, with 5,672, 540 and 629 videos for training, validation and testing, respectively. RWTH-PhoenixT [33] is included for evaluation, which is introduced in “SLT datasets”.

For *SLT*, we make evaluation on RWTH-PhoenixT [33] dataset, which is treated as the extended version of RWTH-Phoenix. It provides parallel sign gloss and translation annotations, to evaluate both continuous SLR and SLT tasks. It contains 8,257 videos, which are divided into three sets: 7,096 for training, 519 for validation, and 642 for testing. RWTH-Phoenix and RWTH-PhoenixT are both recorded at the resolution of 210×260 .

4.1.2 Evaluation Metrics

To evaluate whether our framework works during pre-training, we adopt the metrics for evaluating pose estimation accuracy. Specifically, we report the Percentage of Correct Keypoints (PCK) score and the Area Under the Curve (AUC) on the PCK threshold ranging from 20 to 40 pixels. PCK defines the keypoint to be correct if the Euclidean distance between this keypoint and ground truth is lower than the threshold. The distance metric is expressed in pixels.

For *isolated SLR*, We utilize the accuracy metrics, including per-instance (P-I) and per-class (P-C) metrics. P-I and P-C denote the average accuracy over all the instances and classes, respectively. Following previous works [5], [13], we report Top-1 and Top-5 P-I and P-C metrics on MSASL

and WLASL. Since each class in SLR500 contains the same number of samples, P-I is equal to P-C and we only report one of them.

For *continuous SLR*, we utilize Word Error Rate (WER) as the evaluation metric. WER is the editing distance, which measures the least operations (substitution, deletion and insertion) to transform the hypothesis to the reference gloss sentence as follows,

$$WER = \frac{N_i + N_d + N_s}{L}, \quad (21)$$

where N_i , N_d , and N_s are the number of operations for insertion, deletion, and substitution, respectively. L denotes the length of the reference sentence.

For *SLT*, we adopt BLEU [80] and ROUGE [81] metrics which are commonly utilized in neural machine translation. BLEU calculates the overlap rate of n -gram between the generated text and the reference text, and n ranges from 1 to 4. ROUGE is a metric based on the recall rate and measures the sentence-level structure similarity. In this work, we refer to the ROUGE-L F1-Score.

4.1.3 Implementation Details

In our experiment, all the models are implemented by PyTorch [82] and trained on NVIDIA RTX 3090. Since sign language datasets contain no available pose annotations, we utilize MMPose [83] to extract 133 full 2D keypoints, consisting of 23 body joints, 68 face joints and 42 hand joints. The hyper-parameters ϵ , λ , w_β and w_δ are set as 0.5, 0.01, 10.0 and 100.0, respectively. During decoding, the beam width is set as 10 and 3 for continuous SLR and SLT, respectively.

During the pre-training stage, the utilized data includes the training data from all aforementioned sign datasets, along with other collected data from [84], [85]. In total, the pre-training data volume is 230,246 videos. The Adam optimizer is adopted to train the framework for 60 epochs with the weight decay set as 0.01. The learning rate warms up over the first 10% of the training process, and then decays linearly from the peak rate ($1e-4$). All the frames are fed into the framework.

For all the downstream tasks, the Adam optimizer is still adopted. The learning rate for isolated SLR, continuous SLR and SLT are $1e-4$, $1e-4$ and $5e-5$, respectively. For continuous SLR and SLT, we follow the setting [9]. Spatial-temporal data augmentation is utilized during training. *Spatially*, following [24], we adopt random moving augmentation to simulate spatial disturbance induced by rotation, translation and scaling factors. *Temporally*, for isolated SLR, we extract 32 frames from the origin video using random and center sampling for training and testing, respectively. While for continuous SLR and SLT, we randomly sample 80% frames during training and utilize all the frames during testing.

4.2 Ablation Study

In this section, we first validate the framework feasibility. Then we perform detailed ablation studies on different components of our framework.

TABLE 1

Framework feasibility validation on HANDS17. “P@20” denotes the PCK metrics with the error threshold set as 20 pixel. “Joint”, “Frame” and “Clip” denote the masked joint modeling, masked frame modeling and masked clip modeling, respectively. “Input” and “Output” represent the corrupted input pose and the reconstructed pose sequence by our framework, respectively.

Joint	Mask		Input		Output	
	Frame	Clip	P@20	AUC	P@20	AUC
✓			90.06	89.99	94.60	95.15
	✓		74.83	74.80	93.30	95.13
		✓	60.99	60.00	91.94	93.47
✓	✓	✓	66.65	66.63	94.00	94.74

TABLE 2

Impact of the Transformer layers N on MSASL dataset. N denotes the number of the Transformer encoder layers in our framework.

N	100		200		1000	
	P-I	P-C	P-I	P-C	P-I	P-C
2	82.56	82.35	74.47	75.51	59.20	56.70
3	84.94	85.23	78.51	79.35	62.42	60.15
4	83.75	83.56	76.97	77.74	60.69	57.34
5	83.88	84.23	77.04	77.93	61.27	58.30

4.2.1 Framework Feasibility

We validate the framework feasibility via observing its pose reconstruction capability, on HANDS17 dataset with hand pose annotation available. In this setting, we adopt all masked modeling strategies to train our framework on this dataset. During validation, we perform different masking cases on the input sequence and evaluate the framework output quality. As shown in Table 1, the output metrics are higher than those of the input under all masking cases, which validates our framework feasibility. Besides, we qualitatively visualize the hand pose reconstruction in Figure 4. It can be observed that the reconstructed hand pose sequence is consistent with the ground truth, even under the severely corrupted input situation. It is largely attributed to inherent contextual cues captured by our framework via our designed pretext task.

4.2.2 Ablation Study

To study the impact of different hyper-parameters and settings in our approach, we conduct experiments on MSASL and its subset with per-instance and per-class Top-1 accuracy as the performance indicator.

Impact of the Transformer layers N . As shown in Table 2, the accuracy gets improved when the number N of Transformer layers increases. It reaches the peak when N is equal to 3. There exists difference in the best N between the original BERT and ours, which may be attributed to different characteristics between the sign pose and text domains. In all the experiments, we set N as 3 unless stated.

Impact of the pose θ dimension in the hand-model-aware decoder. From Table 3, the pose θ dimension represents the MANO characterization capability of the hand gesture. The increase of the pose dimension brings enhanced capability and accuracy improvement on the downstream SLR. It reaches the top when the dimension is equal to 25. However,

TABLE 3

Impact of the pose θ dimension in the hand-model-aware decoder on MSASL dataset.

Dimension	100		200		1000	
	P-I	P-C	P-I	P-C	P-I	P-C
15	82.83	82.83	76.01	76.50	61.65	58.59
25	84.94	85.23	78.51	79.35	62.42	60.15
35	83.88	84.20	77.19	77.99	61.60	59.04

TABLE 4

Impact of the temporal span K in masked clip modeling on MSASL dataset. K represents that the masked clip duration ranges from 2 to K .

K	100		200		1000	
	P-I	P-C	P-I	P-C	P-I	P-C
4	81.90	82.17	73.58	74.17	60.43	57.51
8	83.88	83.55	76.60	77.57	61.94	59.76
12	82.96	82.83	74.98	75.16	60.93	58.67

TABLE 5

Impact of different temporal information extraction on MSASL dataset.

Method	100		200		1000	
	P-I	P-C	P-I	P-C	P-I	P-C
PE	84.94	85.23	78.51	79.35	62.42	60.15
GCN_Tem-3	83.36	83.41	76.09	76.85	60.21	57.62
GCN_Tem-5	83.62	84.36	76.09	77.12	60.55	58.13

further increasing does not bring more performance gain, which may be caused by the optimization difficulty.

Impact of the temporal span K in masked clip modeling. In Table 4, it can be observed that the accuracy reaches the top when K is equal to 8. The suitable temporal mask span enforces the framework to capture the temporal dynamics during sign language. In the following, we set K as 8 unless stated.

Impact of different temporal information extraction on MSASL dataset. There are many alternative methods to extract temporal information. Besides temporal position encoding, directly extracting temporal information is also a solution. To this end, we modify the current GCN into the temporal variant following the practice in [74]. Specifically, the original spatial GCN graph is replaced with the spatial-temporal one via adding the local connections along the temporal dimension. With this modification, the gesture state extractor embeds the temporal receptive field of additional k adjacent input frames and thus captures the temporal information. Besides, since our pretext task needs to recover the masked pose token in the corresponding output, we utilize padding to keep the sequence length after the gesture extractor the same as the input.

As shown in Table 5, we perform comparison on these different temporal extraction methods. “PE” denotes utilizing the position encoding for temporal information extraction. For “GCN_Tem- k ”, we remove the temporal position encoding and directly extract temporal information via our modified GCN backbone. k represents the number of adjacent frames. These settings achieve comparable performance on the downstream SLR. It can be explained that the following Transformer encoder contains the strong capability

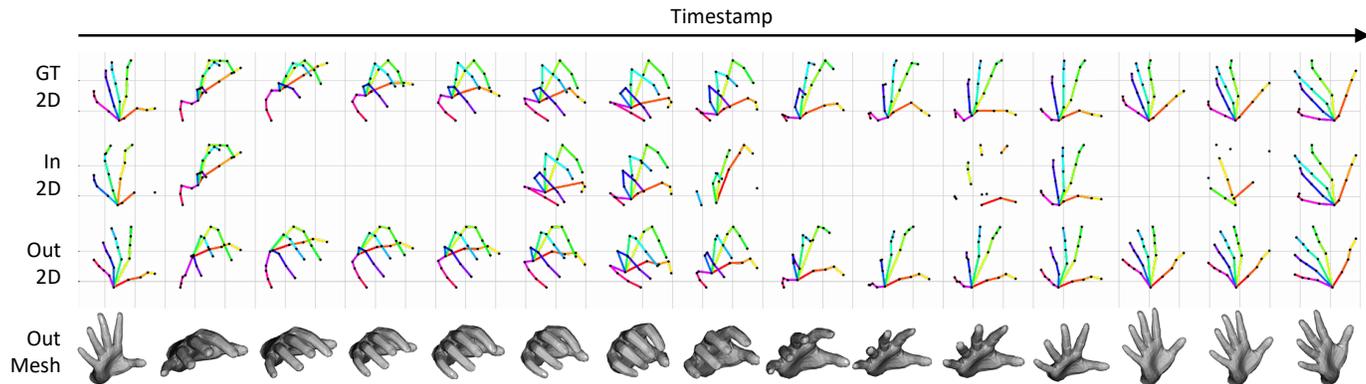


Fig. 4. Qualitative illustration of the framework feasibility on HANDS17. We exhibit 15 continuous frames of a video. Four rows represent the ground truth pose, input pose disturbed by all kinds of masked modeling strategies (joint, frame and clip), reconstructed sequence and the middle mesh representation, respectively. Notably, blanks in the second row represent all joints in the corresponding frames are masked.

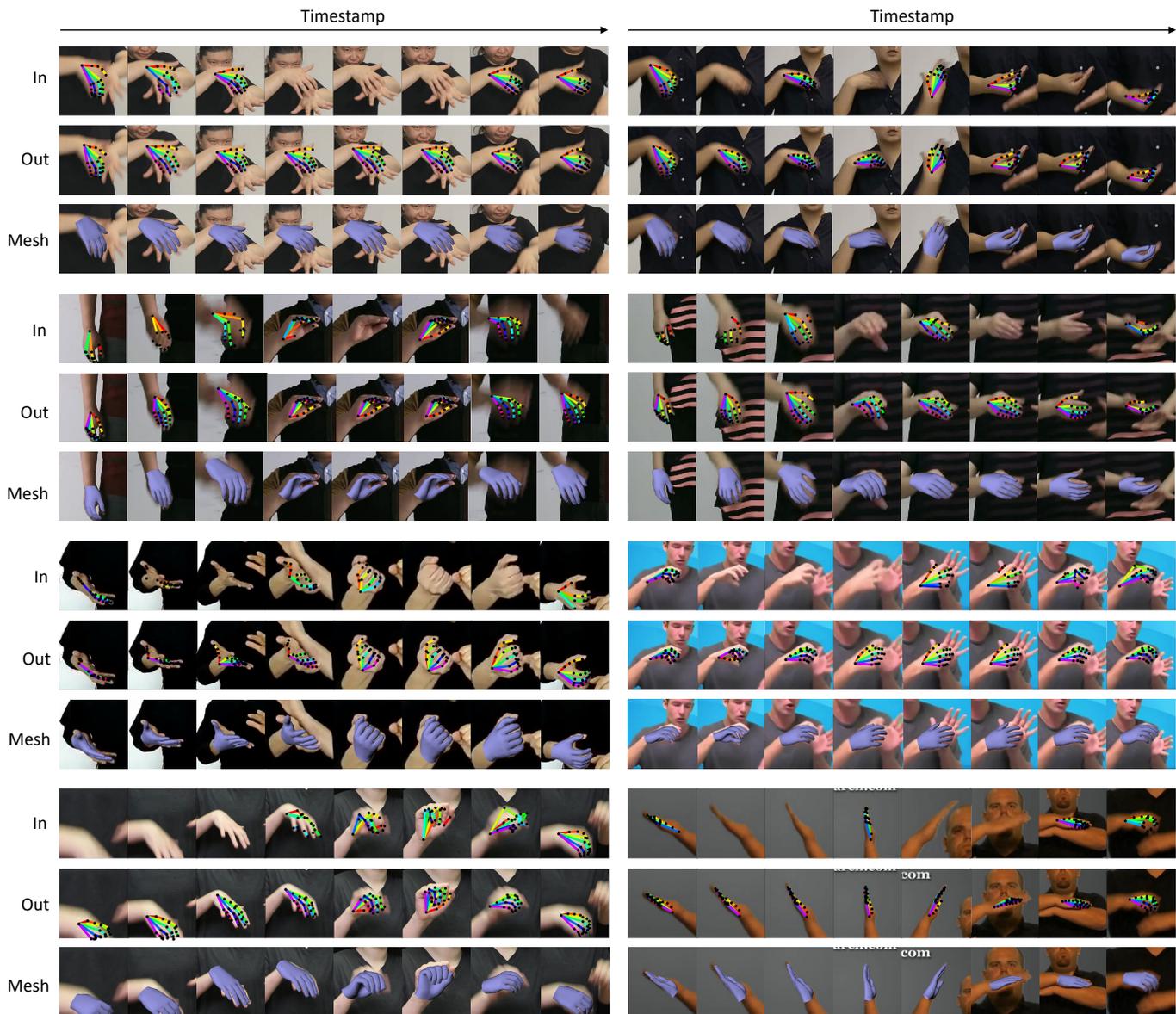


Fig. 5. More visualization samples under sign data sources with no hand pose annotation. For each sample, 8 continuous frames are visualized. “In”, “Out” and “Mesh” denote the input hand pose, the reconstructed hand pose and the intermediate hand mesh, respectively. For clarity, we visualize all the poses and meshes on their aligned RGB image planes.

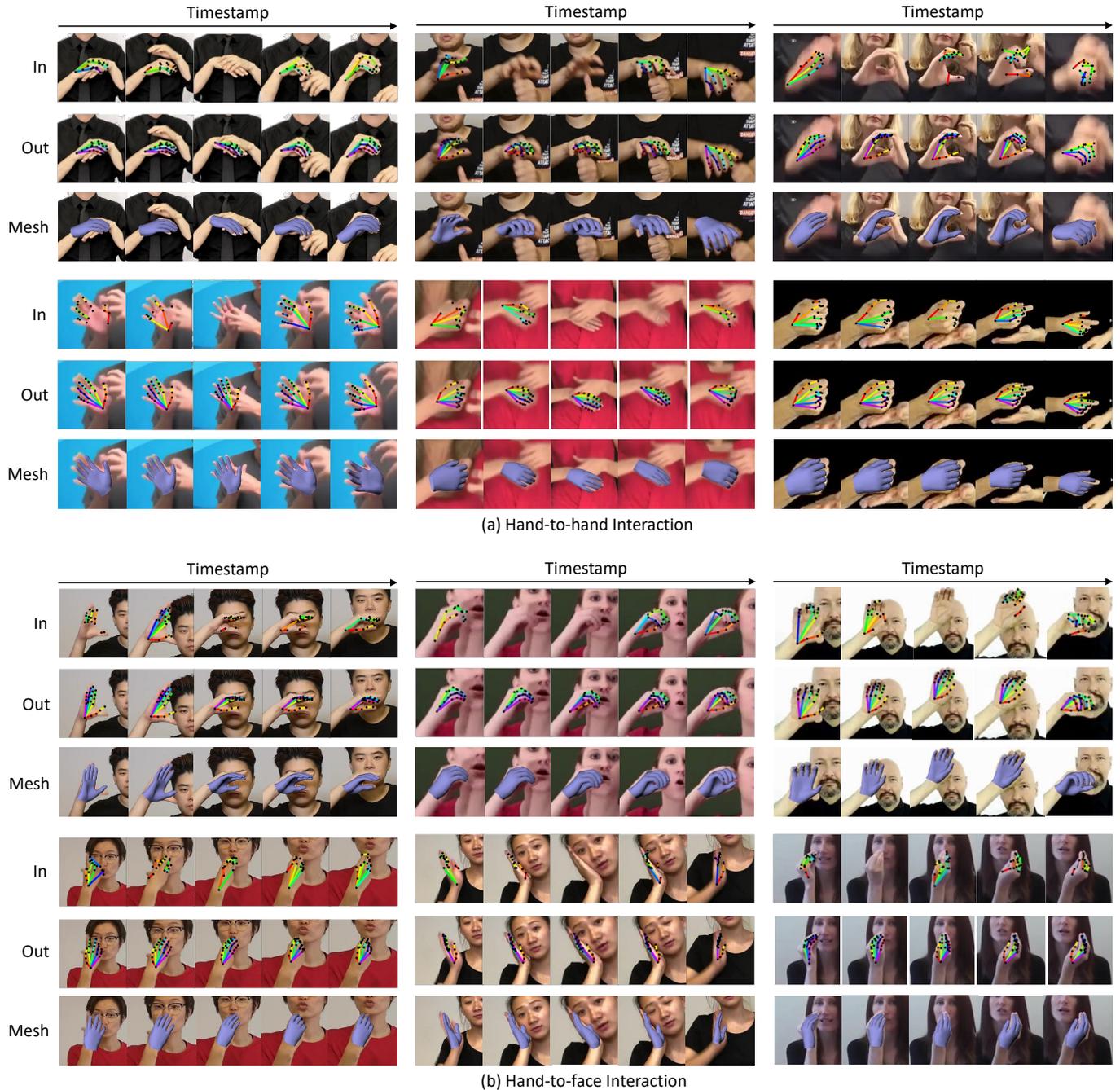


Fig. 6. More visualization samples on two types of hard interaction cases during sign language expression, *i.e.*, hand-to-hand interaction and hand-to-face interaction. For each sample, 5 continuous frames are visualized. “In”, “Out” and “Mesh” denote the input hand pose, the reconstructed hand pose and the intermediate hand mesh, respectively. For clarity, we only plot one hand and visualize its poses and meshes on their aligned RGB image planes.

mediate representation, also improves the interpretability of our method.

As illustrated in Figure 6, we further demonstrate qualitative results on two types of hard cases, *i.e.*, hand-to-hand interaction and hand-to-face interaction. Due to the similar appearance of hand and face and complex self- and mutual occlusion, these interactions bring inherent ambiguity and cause failure in hand pose estimation. Even under these hard cases, our framework can rectify the noisy inputs and infer all the poses which well align the image plane. This strong hallucination capability may be largely attributed to

the well-modeled statistics in the sign language domain.

4.4 Comparison with Other Pre-Training Strategies

As demonstrated in Table 11, we compare with other pre-training strategies, including supervised and state-of-the-art self-supervised pre-training methods. For fair comparison, we pre-train on the same SignBERT+ encoder backbone (embedding layer and Transformer encoder).

Similar to [86], supervised pre-training denotes that we add a classifier (an MLP and softmax layer) on top of the backbone and perform pre-training under the classification

TABLE 11

Comparison with other pre-training strategies on downstream tasks. For fair comparison, all the pre-training methods are performed on the same backbone, *i.e.*, embedding layer and Transformer encoder. The first row represents the framework is directly fine-tuned on the downstream tasks without pre-training. “Partial” and “All” denote the corresponding classification data and all pre-training data, respectively. The data volumes of “Partial” and “All” are about 160k and 230k videos, respectively. (\uparrow denotes the higher the better, while \downarrow represents the lower the better.)

Method	Pre-Train	MSASL				WLASL				SLR500	RWTH-Phoenix		RWTH-PhoenixT	
		P-I		P-C		P-I		P-C		P-I	Dev	Test	Dev	Test
		Top-1 \uparrow	Top-5 \uparrow	Top-1 \uparrow	WER \downarrow	WER \downarrow	WER \downarrow	WER \downarrow						
Baseline	Scratch	53.31	75.98	50.43	74.12	38.33	72.59	36.40	71.23	92.1	43.6	43.4	42.2	42.6
Supervised	Partial	58.82	80.42	56.27	79.54	46.00	79.95	43.63	78.32	94.7	42.5	42.6	40.8	41.7
V-MoCo [86]	Partial	54.22	78.26	51.31	77.03	39.12	72.79	36.93	71.15	94.1	42.1	43.4	40.3	40.8
Ours	Partial	60.13	82.12	57.19	80.72	47.46	81.62	45.03	80.31	95.1	36.7	36.2	35.0	35.2
V-MoCo [86]	All	55.27	79.72	52.06	78.31	40.17	75.19	37.71	73.38	94.8	41.1	41.3	39.2	40.0
Ours	All	62.42	83.49	60.15	82.44	48.85	82.48	46.37	81.33	95.4	34.0	34.1	32.9	33.6

Method	Pre-Train	RWTH-PhoenixT									
		Dev					Test				
		ROUGE \uparrow	BLEU-1 \uparrow	BLEU-2 \uparrow	BLEU-3 \uparrow	BLEU-4 \uparrow	ROUGE \uparrow	BLEU-1 \uparrow	BLEU-2 \uparrow	BLEU-3 \uparrow	BLEU-4 \uparrow
Baseline	Scratch	40.23	39.09	26.40	19.63	15.50	39.83	39.27	26.98	20.10	15.90
Supervised	Partial	41.63	40.87	28.06	20.84	16.49	41.58	41.60	28.82	21.68	17.29
V-MoCo [86]	Partial	42.44	41.94	28.97	21.43	16.84	41.75	41.83	28.82	21.36	16.82
Ours	Partial	44.15	44.03	31.18	23.71	19.00	44.33	43.51	30.92	23.59	19.10
V-MoCo [86]	All	42.79	42.55	29.40	22.13	17.53	41.95	42.68	29.49	21.89	17.36
Ours	All	45.53	44.45	31.88	24.59	19.86	44.89	44.35	32.09	24.92	20.41

task. Specifically, supervised pre-training is conducted on a portion of the original pre-training data (denoted as “Partial”), *i.e.*, the corresponding classification (isolated SLR) benchmarks. The “Partial” and “All” data volumes are 160,113 and 230,246 videos, respectively.

For fair comparison with supervised pre-training, we define two evaluation settings for self-supervised methods, *i.e.*, pre-training on the “Partial” and “All” data volumes. We adopt the state-of-the-art self-supervised pre-training method, *i.e.*, V-MoCo [86]. It is based on contrastive learning, and can be treated as the extended version of MoCo [49] into the video domain. Since it originally works on the RGB domain, we make a few modifications by replacing its backbone with ours, which is able to process the pose modality. During pre-training, we randomly sample two clips with 32 consecutive frames from the same sign pose sequence, as the query and positive samples. The negative samples are obtained from the clips of other videos. During training, its objective is to maximize the similarity between the query and positive samples, which aims to learn temporally persistent features of the same video.

We evaluate the effectiveness of these pre-training methods on all three downstream tasks, *i.e.*, isolated SLR (MSASL, WLASL, SLR500), continuous SLR (RWTH-Phoenix and RWTH-PhoenixT) and SLT (RWTH-PhoenixT). In isolated SLR, supervised pre-training brings a larger performance gain than V-MoCo. While for continuous SLR and SLT, the supervised pre-training method brings a relatively smaller performance gain, whose performance is worse than that of V-MoCo. Supervised pre-training exhibits the limited generalization capability to the downstream tasks, whose objective is inconsistent with classification. For self-supervised pre-training methods, *i.e.*, V-MoCo and Ours, they do not rely on annotated data and scale well with larger pre-training data volume. Notably, when compared with other pre-training strategies, our method achieves the best performance on all downstream tasks with notable gains.

4.5 Comparison with State-of-the-art Methods

In this section, we compare our method with previous state-of-the-art methods on three main downstream tasks, including isolated SLR, continuous SLR and SLT. For comparison, we group them into pose-based and RGB-based methods.

4.5.1 Isolated Sign Language Recognition

Evaluation on MSASL [17]. MSASL introduces new challenges given its unconstrained recording conditions. As illustrated in Table 12, the accuracy of previous pose-based methods lags largely behind the RGB-based counterpart. It is mainly attributed to the pose detection failure caused by partially occluded upper body, motion blur and complex background, *etc.* TCK [7] and BSL [13] propose different pre-training techniques on the I3D backbone by leveraging external sign data. Our method achieves new state-of-the-art performance under both pose-based and RGB-based comparison settings with a notable gain. Notably, when compared with previous SignBERT [5], our framework outperforms it by 12.88% per-instance Top-1 accuracy improvement on the full set.

Evaluation on WLASL [18]. WLASL is also the unconstrained recording setting. Compared with MSASL, it is more challenging with fewer samples and doubled vocabulary size. As shown in Table 13, it is worth mentioning that our single pose-based method even outperforms the most challenging RGB-based method [7], [13], with over 2% Top-1 per-instance accuracy improvement on all sets. When fused with the RGB baseline, the performance of our method further gets improved.

Evaluation on SLR500 [16]. As demonstrated in Table 14, STIP [87] and GMM-HMM [88] are the traditional methods based on hand-crafted features. Since this dataset is recorded under the controlled setting, the performance is quite saturated with only Top-1 accuracy reported. GLE-Net [84] is a challenging method, which performs feature enhancement from the global and local views. Our method achieves 97.8% Top-1 accuracy, which is new state-of-the-art performance.

TABLE 12
Evaluation of isolated SLR on MSASL dataset (the higher the better). [17] denotes the RGB baseline for fusion.

Method	MSASL100				MSASL200				MSASL			
	Per-instance		Per-class		Per-instance		Per-class		Per-instance		Per-class	
	Top-1	Top-5										
Pose-based												
ST-GCN [24]	59.84	82.03	60.79	82.96	52.91	76.67	54.20	77.62	36.03	59.92	32.32	57.15
SignBERT [5]	76.09	92.87	76.65	93.06	70.64	89.55	70.92	90.00	49.54	74.11	46.39	72.65
Ours	84.94	95.77	85.23	95.76	78.51	92.49	79.35	93.03	62.42	83.49	60.15	82.44
RGB-based												
I3D [17]	-	-	81.76	95.16	-	-	81.97	93.79	-	-	57.69	81.05
HMA [72]	73.45	89.70	74.59	89.70	66.30	84.03	67.47	84.03	49.16	69.75	46.27	68.60
TCK [7]	83.04	93.46	83.91	93.52	80.31	91.82	81.14	92.24	-	-	-	-
BSL [13]	-	-	-	-	-	-	-	-	64.71	85.59	61.55	84.43
SignBERT (+ R) [5]	89.56	97.36	89.96	97.51	86.98	96.39	87.62	96.43	71.24	89.12	67.96	88.40
Ours (+ R) [5]	90.75	97.75	91.52	97.73	88.08	96.47	88.62	96.47	73.71	90.12	70.77	89.30

TABLE 13
Evaluation of isolated SLR on WLASL dataset (the higher the better). I3D [18] denotes the RGB baseline for fusion.

Method	WLASL100				WLASL300				WLASL			
	Per-instance		Per-class		Per-instance		Per-class		Per-instance		Per-class	
	Top-1	Top-5										
Pose-based												
ST-GCN [24]	50.78	79.07	51.62	79.47	44.46	73.05	45.29	73.16	34.40	66.57	32.53	65.45
Pose-TGCN [18]	55.43	78.68	-	-	38.32	67.51	-	-	23.65	51.75	-	-
PSLR [23]	60.15	83.98	-	-	42.18	71.71	-	-	-	-	-	-
SignBERT [5]	76.36	91.09	77.68	91.67	62.72	85.18	63.43	85.71	39.40	73.35	36.74	72.38
Ours	79.84	92.64	80.72	93.08	73.20	90.42	73.77	90.58	48.85	82.48	46.37	81.33
RGB-based												
I3D [18]	65.89	84.11	67.01	84.58	56.14	79.94	56.24	78.38	32.48	57.31	-	-
HMA [72]	-	-	-	-	-	-	-	-	37.91	71.26	35.90	70.00
TCK [7]	77.52	91.08	77.55	91.42	68.56	89.52	68.75	89.41	-	-	-	-
BSL [13]	-	-	-	-	-	-	-	-	46.82	79.36	44.72	78.47
SignBERT (+ R) [5]	82.56	94.96	83.30	95.00	74.40	91.32	75.27	91.72	54.69	87.49	52.08	86.93
Ours (+ R)	84.11	96.51	85.05	96.83	78.44	94.31	79.12	94.43	55.59	89.37	53.33	88.82

TABLE 14
Evaluation of isolated SLR on SLR500 dataset (the higher the better). [84] denotes the RGB baseline for fusion.

Method	Accuracy
Pose-based	
ST-GCN [24]	90.0
SignBERT [5]	94.5
Ours	95.4
RGB-based	
STIP [87]	61.8
GMM-HMM [88]	56.3
3D-R50 [89]	95.1
HMA [72]	95.9
GLE-Net [84]	96.8
SignBERT (+ R) [5]	97.7
Ours (+ R)	97.8

In summary, our method greatly shrinks the performance gap between pose-based and RGB-based methods on isolated SLR. Under the challenging in-the-wild conditions, our method even outperforms the challenging RGB-based methods. It can be attributed to our designed masking modeling strategies and incorporated prior during pre-training.

4.5.2 Continuous Sign Language Recognition

Evaluation on RWTH-Phoenix [14]. As demonstrated in Table 16, we exhibit experiment results on RWTH-Phoenix dataset. Due to the lack of pose-based methods, we adopt two representative RGB-based methods [9], [29] by only changing its visual encoder with the GCN to process pose modality, denoted as “P-BLSTM” and “P-Trans”. Pose-based methods lag largely behind RGB-based methods, which is largely caused by pose failure caused by the low-quality data and motion blur. Among pose-based methods, our method largely outperforms the most challenging competitor P-BLSTM with 5.6% and 5.2% WER improvement on the dev and test set, respectively. When fused with the RGB baseline, our method achieves new state-of-the-art performance, *i.e.*, 19.9% and 20.0% WER on the dev and test set, respectively.

Evaluation on RWTH-PhoenixT [33]. We make comparison on RWTH-PhoenixT in Table 17. This dataset additionally provides spoken German translation corresponding to the sign gloss annotation. [15] utilizes the spoken translation to infer the mouth shape label, which provides auxiliary cues to recognition. Besides, it releases multi-stream versions for further performance improvement. STMC [35] also leverages the multi-cue information from the full frame, hand,

TABLE 15
Evaluation of SLT on RWTH-PhoenixT dataset (the higher the better). [36] denotes the RGB baseline method for fusion.

Method	Dev					Test				
	ROUGE	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE	BLEU-1	BLEU-2	BLEU-3	BLEU-4
Pose-based										
Skeletor [90]	32.66	31.97	19.53	14.01	10.91	31.80	31.86	19.11	13.49	10.35
Ours	45.53	44.45	31.88	24.59	19.86	44.89	44.35	32.09	24.92	20.41
RGB-based										
Sign2Text [33]	31.80	31.87	19.11	13.16	9.94	31.80	32.24	19.03	12.83	9.58
TSPNet [34]	-	-	-	-	-	34.96	36.10	23.12	16.88	13.41
MCT [91]	45.90	-	-	-	19.51	43.57	-	-	-	18.51
SL-Trans [9]	-	47.26	34.40	27.05	22.38	-	46.61	33.73	26.19	21.32
BN-TIN-Trans [36]	46.87	46.90	33.98	26.49	21.78	46.98	47.57	34.64	26.78	21.68
SimulSLT [92]	36.04	36.01	22.60	16.05	12.39	35.13	35.92	22.70	16.03	12.27
PiSLTRc-T [92]	47.89	46.51	33.78	26.78	21.48	48.13	46.22	33.56	26.04	21.29
STMC [35]	48.24	47.60	36.43	29.18	24.08	46.65	46.98	36.09	28.70	23.65
SignBT [36]	50.29	51.11	37.90	29.80	24.45	49.54	50.80	37.75	29.72	24.32
Ours (+ R)	51.12	51.46	38.28	30.30	24.95	50.63	52.01	39.19	31.06	25.70

TABLE 16
Evaluation of continuous SLR on RWTH-Phoenix dataset (the lower the better). [29] denotes the RGB baseline for fusion.

Methods	Dev			Test		
	del / ins	WER		del / ins	WER	
Pose-based						
P-BLSTM [29]	13.4 / 3.5	39.6		12.3 / 3.4	39.3	
P-Trans [9]	16.0 / 3.2	40.9		14.9 / 3.4	40.4	
Ours	9.0 / 6.3	34.0		7.9 / 6.0	34.1	
RGB-based						
CMLLR [14]	21.8 / 3.9	55.0		20.3 / 4.5	53.0	
1-Million-Hand [93]	16.3 / 4.6	47.1		15.2 / 4.6	45.1	
CNN-Hybrid [94]	12.6 / 5.1	38.3		11.1 / 5.7	38.8	
SubUNets [76]	14.6 / 4.0	40.8		14.3 / 4.0	40.7	
RCNN [95]	13.7 / 7.3	39.4		12.2 / 7.5	38.7	
Re-Sign [25]	-	27.1		-	26.8	
Hybrid CNN-HMM [26]	-	31.6		-	32.5	
CNN-LSTM-HMM [15]	-	26.0		-	26.0	
CTF [96]	12.8 / 5.2	37.9		11.9 / 5.6	37.8	
Dilated [97]	8.3 / 4.8	38.0		7.6 / 4.8	37.3	
IAN [98]	12.9 / 2.6	37.1		13.0 / 2.5	36.7	
DNF [29]	7.8 / 3.5	23.8		7.8 / 3.4	24.4	
FCN [99]	-	23.7		-	23.9	
CMA [30]	7.3 / 2.7	21.3		7.3 / 2.4	21.9	
PiSLTRc-R [100]	8.1 / 3.4	23.4		7.6 / 3.3	23.2	
STMC [35]	7.7 / 2.4	21.7		7.4 / 2.6	20.7	
VAC [8]	7.9 / 2.5	21.2		8.4 / 2.6	22.3	
Ours (+ R)	4.8 / 3.7	19.9		4.2 / 3.8	20.0	

face and pose and becomes the most challenging competitor. Our method (Ours (+ R)) outperforms it while only utilizing the full video and pose information.

4.5.3 Sign Language Translation

Evaluation on RWTH-PhoenixT [33]. As shown in Table 15, we perform comparison on RWTH-PhoenixT dataset, which is the current most popular benchmark for evaluating SLT. Contemporaneous with ours, Skeletor [90] is an influential work that conducts BERT style pre-training but in another field of pose estimation. Specifically, it inflates the detected poses to 3D ones and conducts BERT-style pre-training with the aim of refining 3D poses. Then it validates its effectiveness on downstream SLT with the refined poses as input. Compared with this challenging pose-based method,

TABLE 17
Evaluation of continuous SLR on RWTH-PhoenixT dataset (the lower the better). [29] denotes the RGB baseline for fusion.

Methods	Dev		Test	
	del / ins	WER	del / ins	WER
Pose-based				
P-BLSTM [29]	13.8 / 3.3	40.2	12.9 / 3.1	40.2
P-Trans [9]	12.9 / 3.7	39.4	11.4 / 3.8	39.8
Ours	9.2 / 4.9	32.9	8.4 / 5.3	33.6
RGB-based				
1-stream [15]	-	24.5	-	26.5
3-stream [15]	-	22.1	-	24.1
DNF [29]	10.5 / 1.9	22.7	9.8 / 2.4	23.5
SL-Trans [9]	11.7 / 6.5	24.9	11.2 / 6.1	24.6
FCN [99]	-	23.3	-	25.1
PiSLTRc-R [100]	4.9 / 4.2	21.8	5.1 / 4.4	22.9
STMC [35]	-	19.6	-	21.0
Ours (+ R)	4.8 / 3.3	18.8	4.3 / 3.9	19.9

our framework directly models the SL statics in the latent semantics space and surpasses it with a larger performance gain, *i.e.*, 8.95% and 10.06% BLEU-4 improvement on the dev and test set, respectively. When compared with RGB-based methods, Ours (+ R) also achieves new state-of-the-art performance, achieving 24.95% and 25.70% BLEU-4 on the dev and test set, respectively.

4.6 Evaluation with Deaf Community

The end goal of automatic sign language understanding is to make the daily life of the deaf community more convenient. Xu *et al.* [101] make the first attempt of user study to evaluate the built sign gloss dictionary for sign language learners. Evaluation participation of the deaf community is also crucial to better analyze our method and outline future work. In our work, the effectiveness of pre-training is evaluated on the downstream tasks and it is also desirable to perform more direct evaluation on pre-training.

To this end, we conduct a user study with the Institutional Review Board (IRB) approval from our college with granted number No.202200603. This user study aims to analyze the robustness of our framework under different input

TABLE 18

User study with deaf community on robustness of our pre-training model under different input noise levels. The rate represents the average correct rate which means the deaf volunteer is able to correctly identify the semantic meaning via observing the output pose of our framework.

Noise Intensity	Correct Rate	
	Input	Output
0.1	23%	92%
0.2	16%	85%
0.3	10%	79%
0.4	4%	75%

noise levels via evaluating the semantic preservation of the pre-training framework output (pose sequence). Different noise levels are achieved via choosing different portions of the input tokens to add noise. There are 10 deaf volunteers participating in this study. In the study, each volunteer is asked to judge whether the corresponding sign gloss can be identified via observing the framework output. We report the correct rate to indicate semantic preservation. Totally, 100 real-world sign videos are collected and involved in this study. For each video, there are 8 samples, *i.e.*, 4 (noise levels) \times 2 (input and output) needed to be evaluated.

As shown in Table 18, it can be observed that the semantics of the output are well-preserved when the input noise intensity increases, which validates the robustness of our framework. Meanwhile, the correct rate of the output is consistently better than that of the input under all noise levels. It reveals that the modeled statistics via our pre-training can bring positive gains on the semantics. Besides, these deaf participants also give us feedback on the reasons of failed recognition from their perspectives, *e.g.* pose jittering, nonstandard gesture, *etc.* These results can give us some hints on further improving the pre-training design, *e.g.* inserting the basic gesture types of sign language as the constraint.

4.7 Analysis & Future Work

The core of our work is modeling the statistics in the sign language domain via maximizing the likelihood of the joint probability distribution, which benefits the downstream sign language understanding tasks. Despite the success of BERT in NLP, it is non-trivial to leverage its masked modeling pretext task into sign language understanding due to different characteristics between these two domains. Among them, the major one is information density [60]. Originally, the languages are highly semantic and well represented with 1D sequences of text words, which are defined with clarified semantics. In contrast, sign pose, expressed in 3D continuous coordinates, is a kind of well-structured data with both spatial and temporal redundancy. Besides, this kind of signal usually contains noise due to failure estimation. This fundamental difference raises the following issues to resolve, which outlines potential future work.

- Token embedding & Position encoding. These embeddings are needed to be carefully designed considering the hand pose characteristics in the sign language domain, *e.g.* how to effectively represent spatial-temporal positions of hands.

- Masking strategy. It aims to capture the hierarchical context in the sign data, which needs to consider the characteristic of sign pose data.
- Decoder design & Pre-training objective. The decoder in the pre-training stage performs mapping from the latent feature back to the input. In NLP, the decoder predicts the masked discrete words with the cross-entropy objective. While for this task, the goal is to reconstruct the continuous sign hand pose sequence. The involved pre-training objective and decoder are needed to design.

More discussion. In this work, we provide our solution to the above issues and validate the effectiveness of our framework. Pre-training on pose has its pros and cons. Pose data is a semantic and compact representation, which is robust to appearance or background changes and brings potential computation efficiency. On the other hand, our adopted pose input is estimated by the off-the-shelf pose detector. Although our framework embeds the capability to capture the cues from the corrupted input pose sequence, its bottleneck is somewhat limited by the quality of the detected pose. Jointly optimizing the pose detector with our framework may be a possible solution. It is also desirable to extend masked-modeling-based self-supervised pre-training to RGB data. Besides, pre-training can go beyond self-supervised learning, *e.g.* multilingual information may be merged as an auxiliary indicator for better performance on downstream tasks.

5 CONCLUSION

In this paper, we propose the *first* self-supervised pre-trainable framework with hand prior incorporated, namely SignBERT+. Given the dominant role of hand during sign language, we take both hands as visual tokens and carefully embed each visual tokens with gesture state and spatial-temporal position information. Our framework first performs pre-training on a large volume of sign data via reconstructing the masked tokens from the corrupted input sequence. Specifically, we subtly design hierarchical masked modeling strategies (joint, frame and clip). These strategies explicitly consider hand pose characteristics to capture multi-level contextual information. Furthermore, we design the hand-model-aware decoder to incorporate prior for better optimization and context modeling. Then, the pre-trained SignBERT+ is fine-tuned for downstream tasks. Given the task diversities, we design simple yet effective prediction heads on top of the SignBERT+ encoder during fine-tuning. Extensive experiments are conducted among three main video-based sign language understanding tasks, *i.e.* isolated SLR, continuous SLR and SLT. Our experiment results demonstrate the effectiveness of our method, achieving new state-of-the-art performance with a notable gain.

Broader Impact. It is estimated by World Health Organization (WHO) that by 2050 over 700 million people will have disabling hearing loss, which accounts for 10% of global population [102]. The community with hearing loss may feel isolated, lonely and other mental issues when they face the communication barrier in daily life. One way to assist them is to bridge this gap via the automatic sign language understanding technique. Our framework is able to promote

its development. However, our technique is not intended for the potential privacy issue, such as surveillance on sign language communication.

REFERENCES

- [1] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," *arxiv*, pp. 1–12, 2018.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional Transformers for language understanding," in *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2018, pp. 4171–4186.
- [3] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. Salakhutdinov, and Q. V. Le, "XLNet: Generalized autoregressive pretraining for language understanding," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 1–18.
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5999–6009.
- [5] H. Hu, W. Zhao, W. Zhou, Y. Wang, and H. Li, "SignBERT: Pre-training of hand-model-aware representation for sign language recognition," in *International Conference on Computer Vision (ICCV)*, 2021, pp. 11 087–11 096.
- [6] O. Koller, "Quantitative survey of the state of the art in sign language recognition," *arXiv*, pp. 1–40, 2020.
- [7] D. Li, C. Rodriguez, X. Yu, and H. Li, "Transferring cross-domain knowledge for video sign language recognition," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 6205–6214.
- [8] Y. Min, A. Hao, X. Chai, and X. Chen, "Visual alignment constraint for continuous sign language recognition," in *International Conference on Computer Vision (ICCV)*, 2021, pp. 11 542–11 551.
- [9] N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, "Sign language Transformers: Joint end-to-end sign language recognition and translation," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 10 023–10 033.
- [10] T. Pfister, J. Charles, and A. Zisserman, "Large-scale learning of sign language by watching TV (using co-occurrences)." in *British Machine Vision Conference (BMVC)*, 2013, pp. 1–11.
- [11] J. F. Lichtenauer, E. A. Hendriks, and M. J. Reinders, "Sign language recognition by combining statistical DTW and independent classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 30, no. 11, pp. 2040–2046, 2008.
- [12] Y. C. Bilge, R. G. Cinbis, and N. Ikingler-Cinbis, "Towards zero-shot sign language recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, pp. 1–16, 2022.
- [13] S. Albanie, G. Varol, L. Momeni, T. Afouras, J. S. Chung, N. Fox, and A. Zisserman, "BSL-1k: Scaling up co-articulated sign language recognition using mouthing cues," in *European Conference on Computer Vision (ECCV)*, 2020, pp. 35–53.
- [14] 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 (CVIU)*, vol. 141, pp. 108–125, 2015.
- [15] O. Koller, 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 (TPAMI)*, vol. 42, no. 9, pp. 2306–2320, 2020.
- [16] J. Huang, W. Zhou, H. Li, and W. Li, "Attention based 3D-CNNs for large-vocabulary sign language recognition," *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, vol. 29, no. 9, pp. 2822–2832, 2019.
- [17] H. R. V. Joze and O. Koller, "MS-ASL: A large-scale data set and benchmark for understanding american sign language," *British Machine Vision Conference (BMVC)*, pp. 1–16, 2019.
- [18] D. Li, C. Rodriguez, X. Yu, and H. Li, "Word-level deep sign language recognition from video: A new large-scale dataset and methods comparison," in *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020, pp. 1459–1469.
- [19] C. Li, Q. Zhong, D. Xie, and S. Pu, "Co-occurrence feature learning from skeleton data for action recognition and detection with hierarchical aggregation," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2018, pp. 786–792.
- [20] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1110–1118.
- [21] Y. Min, Y. Zhang, X. Chai, and X. Chen, "An efficient PointLSTM for point clouds based gesture recognition," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 5761–5770.
- [22] S. Song, C. Lan, J. Xing, W. Zeng, and J. Liu, "An end-to-end spatio-temporal attention model for human action recognition from skeleton data," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2017, pp. 4263–4270.
- [23] A. Tunga, S. V. Nuthalapati, and J. Wachs, "Pose-based sign language recognition using GCN and BERT," in *WACV Workshop*, 2020, pp. 31–40.
- [24] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2018, pp. 7444–7452.
- [25] O. Koller, S. Zargaran, and H. Ney, "Re-sign: Re-aligned end-to-end sequence modelling with deep recurrent CNN-HMMs," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4297–4305.
- [26] 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 (IJCV)*, vol. 126, no. 12, pp. 1311–1325, 2018.
- [27] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [28] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," in *Workshop on Syntax, Semantics and Structure in Statistical Translation (SSST)*, 2014, pp. 103–111.
- [29] R. Cui, H. Liu, and C. Zhang, "A deep neural framework for continuous sign language recognition by iterative training," *IEEE Transactions on Multimedia (TMM)*, vol. 21, no. 7, pp. 1880–1891, 2019.
- [30] J. Pu, W. Zhou, H. Hu, and H. Li, "Boosting continuous sign language recognition via cross modality augmentation," in *ACM International Conference on Multimedia (ACM MM)*, 2020, pp. 1497–1505.
- [31] B. Shi, A. M. D. Rio, J. Keane, D. Brentari, G. Shakhnarovich, and K. Livescu, "Fingerspelling recognition in the wild with iterative visual attention," in *International Conference on Computer Vision (ICCV)*, 2019, pp. 5400–5409.
- [32] Z. Zhou, V. W. L. Tam, and E. Y. Lam, "Signbert: A bert-based deep learning framework for continuous sign language recognition," *IEEE Access*, vol. 9, pp. 161 669–161 682, 2021.
- [33] N. Cihan Camgoz, S. Hadfield, O. Koller, H. Ney, and R. Bowden, "Neural sign language translation," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7784–7793.
- [34] D. Li, C. Xu, X. Yu, K. Zhang, B. Swift, H. Suominen, and H. Li, "TSPNet: Hierarchical feature learning via temporal semantic pyramid for sign language translation," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 1–12.
- [35] H. Zhou, W. Zhou, Y. Zhou, and H. Li, "Spatial-temporal multi-tuple network for sign language recognition and translation," *IEEE Transactions on Multimedia (TMM)*, vol. 24, pp. 768–779, 2021.
- [36] H. Zhou, W. Zhou, W. Qi, J. Pu, and H. Li, "Improving sign language translation with monolingual data by sign back-translation," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 1316–1325.
- [37] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.
- [38] J. Carreira and A. Zisserman, "Quo vadis, action recognition? A new model and the Kinetics dataset," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6299–6308.
- [39] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, and J. Tang, "Self-supervised learning: Generative or contrastive,"

- IEEE Transactions on Knowledge and Data Engineering (TKDE)*, pp. 1–20, 2021.
- [40] L. Jing and Y. Tian, “Self-supervised visual feature learning with deep neural networks: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 43, no. 11, pp. 4037–4058, 2020.
- [41] R. Zhang, P. Isola, and A. A. Efros, “Colorful image colorization,” in *European Conference on Computer Vision (ECCV)*, 2016, pp. 649–666.
- [42] J. Wang, J. Jiao, L. Bao, S. He, W. Liu, and Y.-H. Liu, “Self-supervised video representation learning by uncovering spatio-temporal statistics,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, pp. 1–16, 2021.
- [43] S. Gidaris, P. Singh, and N. Komodakis, “Unsupervised representation learning by predicting image rotations,” in *International Conference on Learning Representations (ICLR)*, 2018, pp. 1–16.
- [44] M. Noroozi and P. Favaro, “Unsupervised learning of visual representations by solving jigsaw puzzles,” in *European Conference on Computer Vision (ECCV)*, 2016, pp. 69–84.
- [45] K. Wang, L. Lin, C. Jiang, C. Qian, and P. Wei, “3D human pose machines with self-supervised learning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 42, no. 5, pp. 1069–1082, 2019.
- [46] G.-J. Qi, L. Zhang, F. Lin, and X. Wang, “Learning generalized transformation equivariant representations via autoencoding transformations,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 44, no. 4, pp. 2045–2057, 2022.
- [47] I. Misra, C. L. Zitnick, and M. Hebert, “Shuffle and learn: unsupervised learning using temporal order verification,” in *European Conference on Computer Vision (ECCV)*, 2016, pp. 527–544.
- [48] D. Xu, J. Xiao, Z. Zhao, J. Shao, D. Xie, and Y. Zhuang, “Self-supervised spatiotemporal learning via video clip order prediction,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10334–10343.
- [49] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 9729–9738.
- [50] L. Linguo, W. Minsi, N. Bingbing, W. Hang, Y. Jiancheng, and Z. Wenjun, “3D human action representation learning via cross-view consistency pursuit,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [51] X. Chen and K. He, “Exploring simple siamese representation learning,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 15750–15758.
- [52] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International Conference on Machine Learning (ICML)*, 2020, pp. 1597–1607.
- [53] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar *et al.*, “Bootstrap your own latent: A new approach to self-supervised learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 21271–21284.
- [54] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “RoBERTa: A robustly optimized bert pretraining approach,” *arXiv*, pp. 1–13, 2019.
- [55] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, “VideoBERT: A joint model for video and language representation learning,” in *International Conference on Computer Vision (ICCV)*, 2019, pp. 7464–7473.
- [56] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, “VL-BERT: Pre-training of generic visual-linguistic representations,” in *International Conference on Learning Representations (ICLR)*, 2020, pp. 1–16.
- [57] M. Chen, A. Radford, R. Child, J. Wu, H. Jun, D. Luan, and I. Sutskever, “Generative pretraining from pixels,” in *International Conference on Machine Learning (ICML)*, 2020, pp. 1691–1703.
- [58] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations (ICLR)*, 2020, pp. 1–21.
- [59] H. Bao, L. Dong, and F. Wei, “BEiT: Bert pre-training of image transformers,” in *International Conference on Learning Representations (ICLR)*, 2022, pp. 1–16.
- [60] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” *arXiv*, pp. 1–8, 2021.
- [61] S. Sridhar, A. Oulasvirta, and C. Theobalt, “Interactive markerless articulated hand motion tracking using RGB and depth data,” in *International Conference on Computer Vision (ICCV)*, 2013, pp. 2456–2463.
- [62] I. Oikonomidis, M. I. Lourakis, and A. A. Argyros, “Evolutionary quasi-random search for hand articulations tracking,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 3422–3429.
- [63] C. Qian, X. Sun, Y. Wei, X. Tang, and J. Sun, “Realtime and robust hand tracking from depth,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 1106–1113.
- [64] A. Tkach, M. Pauly, and A. Tagliasacchi, “Sphere-meshes for real-time hand modeling and tracking,” *ACM Transactions on Graphics (TOG)*, vol. 35, no. 6, pp. 1–11, 2016.
- [65] I. Oikonomidis, N. Kyriazis, and A. A. Argyros, “Efficient model-based 3D tracking of hand articulations using Kinect,” in *British Machine Vision Conference (BMVC)*, 2011, pp. 1–11.
- [66] L. Ballan, A. Taneja, J. Gall, L. Van Gool, and M. Pollefeys, “Motion capture of hands in action using discriminative salient points,” in *European Conference on Computer Vision (ECCV)*, 2012, pp. 640–653.
- [67] D. Tzionas, L. Ballan, A. Srikantha, P. Aponte, M. Pollefeys, and J. Gall, “Capturing hands in action using discriminative salient points and physics simulation,” *International Journal of Computer Vision (IJCV)*, vol. 118, no. 2, pp. 172–193, 2016.
- [68] J. Romero, D. Tzionas, and M. J. Black, “Embodied hands: Modeling and capturing hands and bodies together,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 6, pp. 1–17, 2017.
- [69] J. P. Lewis, M. Corder, and N. Fong, “Pose space deformation: a unified approach to shape interpolation and skeleton-driven deformation,” in *Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, 2000, pp. 165–172.
- [70] M. Habermann, W. Xu, M. Zollhofer, G. Pons-Moll, and C. Theobalt, “DeepCap: Monocular human performance capture using weak supervision,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 5052–5063.
- [71] A. Boukhayma, R. d. Bem, and P. H. Torr, “3D hand shape and pose from images in the wild,” in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10843–10852.
- [72] H. Hu, W. Zhou, and H. Li, “Hand-model-aware sign language recognition,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2021, pp. 1558–1566.
- [73] G. Moon, S.-I. Yu, H. Wen, T. Shiratori, and K. M. Lee, “Inter-hand2.6m: A dataset and baseline for 3D interacting hand pose estimation from a single RGB image,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 548–564.
- [74] Y. Cai, L. Ge, J. Liu, J. Cai, T.-J. Cham, J. Yuan, and N. M. Thalmann, “Exploiting spatial-temporal relationships for 3D pose estimation via graph convolutional networks,” in *International Conference on Computer Vision (ICCV)*, 2019, pp. 2272–2281.
- [75] L. Kavan and J. Žára, “Spherical blend skinning: a real-time deformation of articulated models,” in *ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games (I3D)*, 2005, pp. 9–16.
- [76] N. C. Camgoz, S. Hadfield, O. Koller, and R. Bowden, “Sub-UNets: End-to-end hand shape and continuous sign language recognition,” in *International Conference on Computer Vision (ICCV)*, 2017, pp. 3075–3084.
- [77] A. Graves, “Generating sequences with recurrent neural networks,” *arXiv*, pp. 1–43, 2013.
- [78] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv*, pp. 1–23, 2016.
- [79] S. Yuan, Q. Ye, G. Garcia-Hernando, and T.-K. Kim, “The 2017 hands in the million challenge on 3D hand pose estimation,” *arXiv*, pp. 1–7, 2017.
- [80] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “BLEU: a method for automatic evaluation of machine translation,” in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2002, pp. 311–318.

- [81] C.-Y. Lin, "ROUGE: A package for automatic evaluation of summaries," in *ACL workshop on Text Summarization Branches Out*, 2004, pp. 74–81.
- [82] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "PyTorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 8026–8037.
- [83] M. Contributors, "OpenMMLab pose estimation toolbox and benchmark," <https://github.com/open-mmlab/mmpose>, 2020.
- [84] H. Hu, W. Zhou, J. Pu, and H. Li, "Global-local enhancement network for NMFs-aware sign language recognition," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 17, no. 3, pp. 1–18, 2021.
- [85] A. Duarte, S. Palaskar, L. Ventura, D. Ghadiyaram, K. DeHaan, F. Metzger, J. Torres, and X. Giro-i Nieto, "How2sign: a large-scale multimodal dataset for continuous american sign language," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 2735–2744.
- [86] C. Feichtenhofer, H. Fan, B. Xiong, R. Girshick, and K. He, "A large-scale study on unsupervised spatiotemporal representation learning," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3299–3309.
- [87] I. Laptev, "On space-time interest points," *International Journal of Computer Vision (IJCV)*, vol. 64, no. 2-3, pp. 107–123, 2005.
- [88] A. Tang, K. Lu, Y. Wang, J. Huang, and H. Li, "A real-time hand posture recognition system using deep neural networks," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 6, no. 2, pp. 1–23, 2015.
- [89] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3D residual networks," in *International Conference on Computer Vision (ICCV)*, 2017, pp. 5533–5541.
- [90] T. Jiang, N. C. Camgoz, and R. Bowden, "Skeleton: Skeletal transformers for robust body-pose estimation," in *CVPR Workshop*, 2021, pp. 3394–3402.
- [91] N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, "Multi-channel transformers for multi-articulatory sign language translation," in *ECCV Workshop*, 2020, pp. 301–319.
- [92] A. Yin, Z. Zhao, J. Liu, W. Jin, M. Zhang, X. Zeng, and X. He, "SimulSLT: End-to-end simultaneous sign language translation," in *ACM International Conference on Multimedia (ACM MM)*, 2021, pp. 4118–4127.
- [93] O. Koller, H. Ney, and R. Bowden, "Deep hand: How to train a cnn on 1 million hand images when your data is continuous and weakly labelled," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 3793–3802.
- [94] O. Koller, O. Zargaran, H. Ney, and R. Bowden, "Deep sign: Hybrid CNN-HMM for continuous sign language recognition," in *British Machine Vision Conference (BMVC)*, 2016, pp. 1–12.
- [95] R. Cui, H. Liu, and C. Zhang, "Recurrent convolutional neural networks for continuous sign language recognition by staged optimization," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 7361–7369.
- [96] S. Wang, D. Guo, W.-g. Zhou, Z.-J. Zha, and M. Wang, "Connectionist temporal fusion for sign language translation," in *ACM International Conference on Multimedia (ACM MM)*, 2018, pp. 1483–1491.
- [97] J. Pu, W. Zhou, and H. Li, "Dilated convolutional network with iterative optimization for continuous sign language recognition," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2018, pp. 885–891.
- [98] —, "Iterative alignment network for continuous sign language recognition," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4165–4174.
- [99] K. L. Cheng, Z. Yang, Q. Chen, and Y.-W. Tai, "Fully convolutional networks for continuous sign language recognition," in *European Conference on Computer Vision (ECCV)*, 2020, pp. 697–714.
- [100] P. Xie, M. Zhao, and X. Hu, "PiSLTRc: Position-informed sign language transformer with content-aware convolution," *IEEE Transactions on Multimedia (TMM)*, pp. 1–13, 2021.
- [101] C. Xu, D. Li, H. Li, H. Suominen, and B. Swift, "Automatic gloss dictionary for sign language learners," in *ACL: System Demonstrations*, 2022, pp. 83–92.
- [102] W. H. Organization. Deafness and hearing loss. (2021, April 1). [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/deafness-and-hearing-loss>