# Multi-Stream Keypoint Attention Network for Sign Language Recognition and Translation

Mo Guan[1], Yan Wang[1*], Guangkun Ma[2], Jiarui Liu[1], Mingzu Sun[1]

[1]School of Information Science and Engineering, Shenyang University of Technology, ShenLiao West Road, Shenyang, 110870, Liaoning, China.
[2]School of Software, Shenyang University of Technology, ShenLiao West Road, Shenyang, 110870, Liaoning, China.

*Corresponding author(s). E-mail(s): wangyan@smail.sut.edu.cn;

**Abstract**

Sign language serves as a non-vocal means of communication, transmitting information and significance through gestures, facial expressions, and bodily movements. The majority of current approaches for sign language recognition (SLR) and translation rely on RGB video inputs, which are vulnerable to fluctuations in the background. Employing a keypoint-based strategy not only mitigates the effects of background alterations but also substantially diminishes the computational demands of the model. Nevertheless, contemporary keypoint-based methodologies fail to fully harness the implicit knowledge embedded in keypoint sequences. To tackle this challenge, our inspiration is derived from the human cognition mechanism, which discerns sign language by analyzing the interplay between gesture configurations and supplementary elements. We propose a multi-stream keypoint attention network to depict a sequence of keypoints produced by a readily available keypoint estimator. In order to facilitate interaction across multiple streams, we investigate diverse methodologies such as keypoint fusion strategies, head fusion, and self-distillation. The resulting framework is denoted as MSKA-SLR, which is expanded into a sign language translation (SLT) model through the straightforward addition of an extra translation network. We carry out comprehensive experiments on well-known benchmarks like Phoenix-2014, Phoenix-2014T, and CSL-Daily to showcase the efficacy of our methodology. Notably, we have attained a novel state-of-the-art performance in the sign language translation task of Phoenix-2014T. The code and models can be accessed at: https://github.com/sutwangyan/MSKA.

**Keywords:** Sign Language Recognition, Sign Language Translation, Self-Attention, Self-Distillation, Keypoint

## 1 Introduction

Sign language, a form of communication utilizing gestures, expressions, and bodily movements, has been the subject of extensive study Bungeroth and Ney (2004); Starner et al (1998); Tamura and Kawasaki (1988). For the deaf and mute community, sign language serves as their primary mode of communication. It holds profound significance, offering an effective medium for this particular demographic to convey thoughts, emotions, and needs, thereby facilitating their active participation in social interactions. Sign language possesses a unique structure, incorporating elements such as the shape, direction, and placement of gestures, along with facial expressions. Its

grammar diverges from that of spoken language, exhibiting differences in grammatical structure and sequence. To address such disparities, certain sign language translation (SLT) tasks integrate gloss sequences before text generation. The transition from visual input to gloss sequences constitutes the process of sign language recognition (SLR). Fig. 1(a) depicts both SLR and SLT tasks.

Gestures play a pivotal role in the recognition and translation of sign language. Indeed, gestures occupy a modest portion of the video, rendering them vulnerable to shifts in the background and swift hand movements during sign language communication. Consequently, this results in challenges in acquiring sign language attributes. Nevertheless, owing to robustness and computational efficiency of gestures, some methodologies advocate for the employment of keypoints to convey it. Ordinarily, sign language videos undergo keypoint extraction using off-the-shelf keypoint estimator. Following this, the keypoint sequences are regionally cropped to be utilized as input for the model, allowing a more precise focus on the characteristics of hand shapes. The TwoStream method, as described in Chen et al (2022b), enhances feature extraction by converting keypoints into heatmaps and implementing 3D convolution. SignBERT+, detailed in the work by Hu et al (2023a), represents hand keypoints as a graphical framework and employs graph convolutional networks for extracting gesture features. Nevertheless, a key drawback of these approaches is the inadequate exploitation of correlation data among keypoints.

To address this challenge, we introduce an innovative network framework that depends entirely on the interplay among keypoints to achieve proficiency in sign language recognition and translation endeavors. Our methodology is influenced by the innate human inclination to prioritize the configuration of gestures and the dynamic interconnection between the hands and other bodily elements in the process of sign language interpretation. The devised multi-stream keypoint attention (MSKA) mechanism is adept at facilitating sign language translation by integrating a supplementary translation network. As a result, the all-encompassing system is designated as MSKA-SLT, as illustrated in Fig. 1(b).

In summary, our contributions primarily consist of the following three aspects:

1. To the best of our knowledge, we are the first to propose a multi-stream keypoint attention, which is built with pure attention modules without manual designs of traversal rules or graph topologies.
2. We propose to decouple the keypoint sequences into four streams, left hand stream, right hand stream, face stream and whole body stream, each focuses on a specific aspect of the skeleton sequence. By fusing different types of features, the model can have a more comprehensive understanding for sign language recognition and translation.
3. We conducted extensive experiments to validate the proposed method, demonstrating encouraging improvements in sign language recognition tasks on the three prevalent benchmarks, *i.e.*, Phoenix-2014 Koller et al (2015), Phoenix-2014T Camgoz et al (2018) and CSL-Daily Zhou et al (2021a). Moreover, we achieved new state-of-the-art performance in the translation task of Phoenix-2014T.

## 2 Related Work

### 2.1 Sign Language Recognition and Translation

Sign language recognition is a prominent research domain in the realm of computer vision, with the goal of deriving sign glosses through the analysis of video or image data. 2D CNNs are frequently utilized architectures in computer vision to analyze image data, and they have garnered extensive use in research pertaining to sign language recognition Cihan Camgoz et al (2017); Niu and Mak (2020); Cui et al (2019); Zhou et al (2021b); Hu et al (2023b,c); Guo et al (2023).

STMC Zhou et al (2021b) proposed a spatio-temporal multi-cue network to address the problem of visual sequence learning. CorrNet Hu et al (2023b) model captures crucial body movement trajectories by analyzing correlation maps between consecutive frames. It employs 2D CNNs to extract image features, followed by a set of 1D CNNs to acquire temporal characteristics. AdaBrowse Hu et al (2023c) introduced a novel adaptive model that dynamically selects the most informative subsequence from the input video sequence by effectively utilizing redundancy modeled for sequential decision tasks. CTCA Guo et al
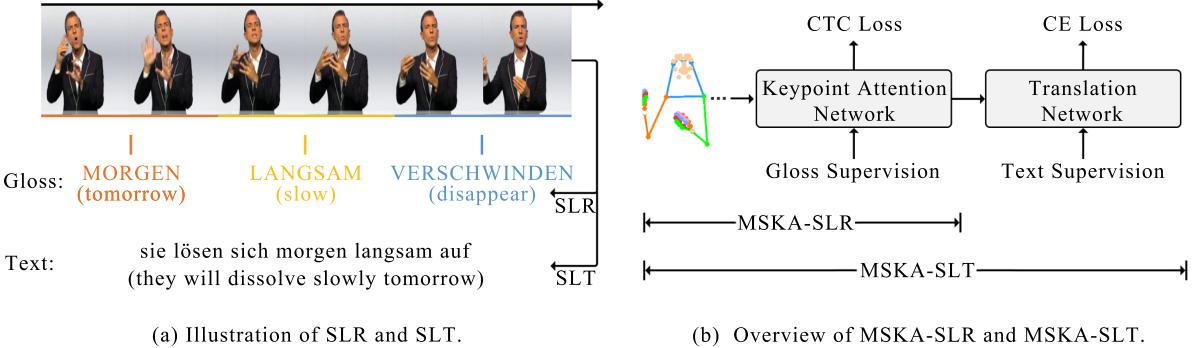
(a) Illustration of SLR and SLT.        (b) Overview of MSKA-SLR and MSKA-SLT.

**Fig. 1**: (a)We choose a sign language video from the Phoenix-2014T dataset and display its gloss sequence alongside the corresponding text. The objective of sign language recognition (SLR) is to instruct models in producing matching gloss representations derived from sign language videos. Conversely, the task of sign language translation (SLT) entails creating textual representations that align with sign language videos. (b) MSKA-SLT is constructed on the foundation of MSKA-SLR to facilitate SLT. Keypoint sequences are depicted in coordinate form.

(2023) build a dual-path network that contains two branches for perceptions of local temporal context and global temporal context. By extending 2D CNNs along the temporal dimension, 3D CNNs can directly process spatio-temporal information in video data. This approach enables a better understanding of the dynamic features of sign language movements, thus enhancing recognition accuracy Li et al (2020); Pu et al (2019); Chen et al (2022a). MMTLB Chen et al (2022a) utilize a pre-trained S3D Xie et al (2018) network to extract features from sign language videos for sign language recognition, followed by the use of a translation network for sign language translation tasks. Recent studies in gloss decoder design have predominantly employed either Hidden Markov Models (HMM) Koller et al (2017, 2018, 2020) or Connectionist Temporal Classification (CTC) Cheng et al (2020); Min et al (2021); Zhou et al (2021b), drawing from their success in automatic speech recognition. We opted for CTC due to its straightforward implementation. While CTC loss offers only modest sentence-level guidance, approaches such as those proposed by Cui et al (2019); Zhou et al (2019); Chen et al (2022b) suggest iteratively deriving detailed pseudo labels from CTC outputs to enhance frame-level supervision. Additionally, Min et al (2021) achieves frame-level knowledge distillation by aligning the entire model with the visual encoder.

In this study, our distillation process leverages the multi-stream architecture to incorporate ensemble knowledge into each individual stream, thereby improving interaction and coherence among the multiple streams. Sign language translation (SLT) involves directly generating textual outputs from sign language videos. Many existing methods frame this task as a neural machine translation (NMT) challenge, employing a visual encoder to extract visual features and feeding them into a translation network for text generation Camgoz et al (2018, 2020b); Chen et al (2022a); Li et al (2020); Zhou et al (2021a); Xie et al (2018); Chen et al (2022b). We adopt mBART Liu et al (2020) as our translation network, given its impressive performance in SLT Chen et al (2022a,b). To attain satisfactory outcomes, gloss supervision is commonly employed in SLT. This involves pre-training the multi-stream attention network on SLR Camgoz et al (2020b); Zhou et al (2021b,a) and jointly training SLR and SLT Zhou et al (2021b,a).

## 2.2 Introduce Keypoints into SLR and SLT

The optimization of keypoints to enhance the efficacy of SLR and SLT remains a challenging issue. Camgoz et al (2020a) introduce an innovative multichannel transformer design. The suggested structure enables the modeling of both inter and
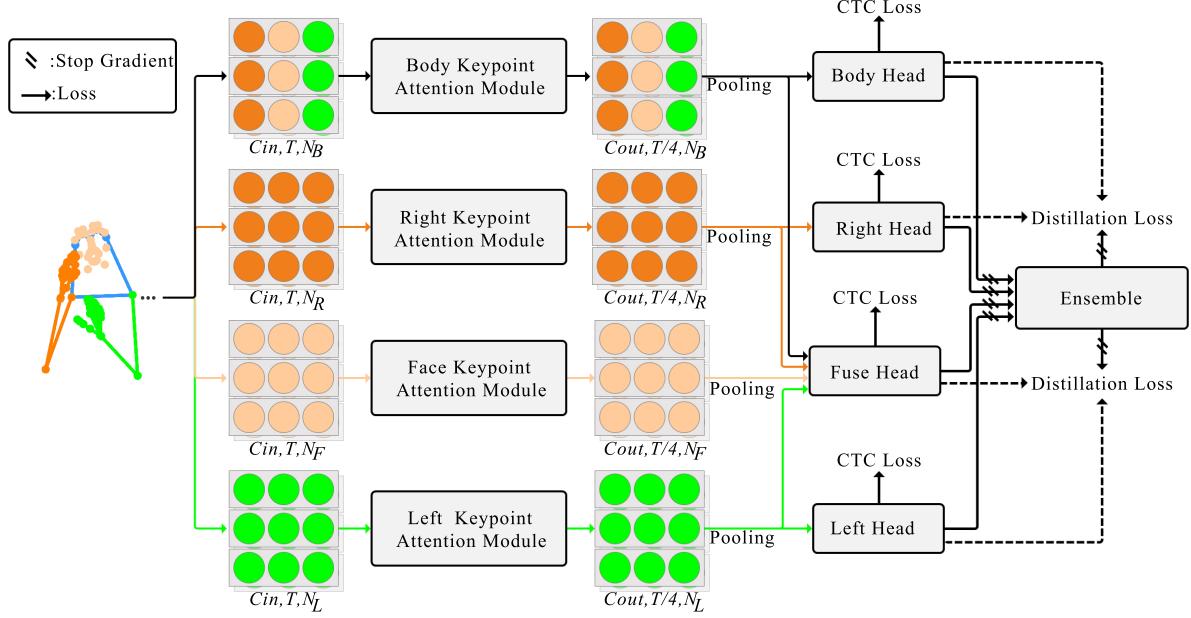
**Fig. 2**: The overview of our MSKA-SLR. The whole network is jointly supervised by the CTC losses and the self-distillation losses. Keypoints are represented in coordinate format.

intra contextual connections among distinct sign articulators within the transformer network, while preserving channel-specific details. Papadimitriou and Potamianos (2020) presenting an end-to-end deep learning methodology that depends on the fusion of multiple spatio-temporal feature streams, as well as a fully convolutional encoder-decoder for prediction. TwoStream-SLR Chen et al (2022b) put forward a dual-stream network framework that integrates domain knowledge such as hand shapes and body movements by modeling the original video and keypoint sequences separately. It utilizes existing keypoint estimators to generate keypoint sequences and explores diverse techniques to facilitate interaction between the two streams. SignBERT+ Hu et al (2023a) incorporates graph convolutional networks (GCN) into hand pose representations and amalgamating them with a self-supervised pre-trained model for hand pose, the aim is to enhance sign language understanding performance. This method utilizes a multi-level masking modeling approach (including joint, frame, and clip levels) to train on extensive sign language data, capturing multi-level contextual information in sign language data. C²SLR Zuo and Mak (2024) aims to ensure coherence between the acquired attention masks

and pose keypoint heatmaps to enable the visual module to concentrate on significant areas.

## 2.3 Self-attention mechanism

Serves as the foundational component within the transformer architecture Vaswani et al (2017); Dai et al (2019), representing a prevalent approach in the realm of natural language processing (NLP). Its operational framework encompasses a set of queries $Q$, keys $K$, and values $V$, each with a dimensionality of $C$, arranged in matrix format to facilitate efficient computation. Initially, the mechanism computes the dot product between the queries and all keys, subsequently normalizing each by $\sqrt{C}$ and applying a softmax function to derive the corresponding weights assigned to the values Vaswani et al (2017). Mathematically, this process can be formulated as follows:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{C}})V \quad (1)$$

## 2.4 Multi-Stream Networks

In this work, our approach directly models keypoint sequences through an attention module. Additionally, to mitigate the issue of data scarcity
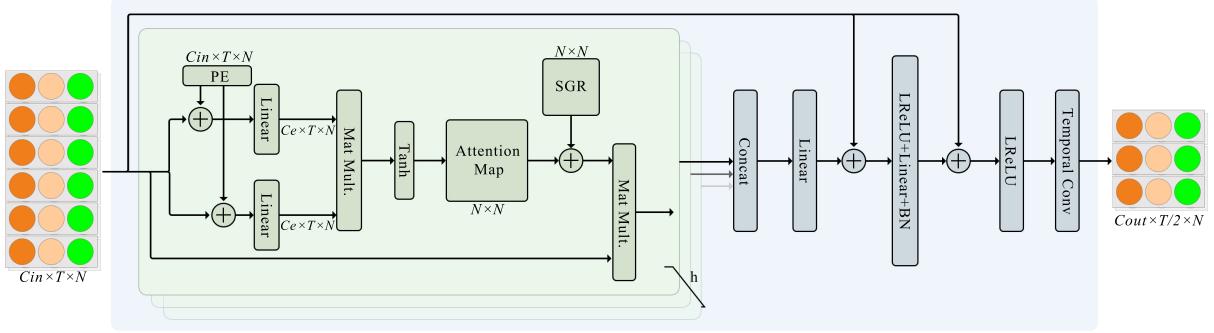
**Fig. 3**: Illustration of the attention module. We show the body attention module as an example. The others attention module is an analogy. The green rounded rectangle box represents a single-head self-attention module. There are totally $h$ self-attention heads, whose output are concatenated and fed into two linear layers to obtain the output. LReLU represents the leaky ReLU Maas et al (2013).

and better capture glosses across different body parts, we introduce multi-stream attention to drive meaningful feature extraction of local features. Modeling the interactions among distinct streams presents a challenging challenge. I3D Carreira and Zisserman (2017) adopts a late fusion strategy by simply averaging the predictions of the two streams. Another approach involves early fusion by lateral connections Feichtenhofer et al (2019), concatenation Zhou et al (2021b), or addition Cui et al (2019) to merge intermediary features of each stream. In this study, we utilize the concept of lateral connections to facilitate mutual supplementation between multiple streams. Additionally, our self-distillation method integrates knowledge from multiple streams into the generated pseudo-targets, thereby achieving a more profound interaction.

# 3 Proposed Method

In this section, we initially present the data augmentation techniques for keypoint sequences. Subsequently, we elaborate on the individual components of MSKA-SLR. Finally, we outline the composition of MSKA-SLT.

## 3.1 Keypoint augment

Typically, sign language video datasets are constrained in size, underscoring the importance of data augmentation. In contrast to prior works such as Guo et al (2023); Hu et al (2023c,b); Chen et al (2022b), our input data comprises

keypoint sequences. Analogous to the augmentation techniques employed in image-related tasks, we implement a step for keypoints: Utilizing HRNet Wang et al (2020) to extract keypoints from sign language videos, wherein the keypoint coordinates are denoted with respect to the top-left corner of the image, with the positive $X$ and $Y$ axes oriented towards the rightward and downward directions, respectively. To utilize data augmentation, we pull the origin back to the center of the image and normalize it by a function: $((x/W, (H-y)/H) - 0.5)/0.5$ , with horizontal to the right and vertical upwards defining the positive directions of the $X$ and $Y$ axes, respectively. Within this context, the variables $x$ and $y$ denote the coordinates of a given point, whereas $H$ and $W$ symbolize the height and width of the image, respectively.

1) We adjust the temporal length of the keypoint sequences within the interval $[\times 0.5\text{-}\times 1.5]$, selecting valid frames randomly from this range. 2) The scaling process involves multiplying the coordinates of each point in the provided keypoint set by a scaling factor. 3) The transformation operation is implemented by applying the provided translation vector to the coordinates of each point in the provided set of keypoint coordinates. 4) During the process of rotation, we achieve this by creating a matrix representing the rotation angle. Given a point $P(x, y)$ in two dimensions, the formula for calculating for the resulting point $P'(x', y')$ with the center at the origin, and a counterclockwise rotation by an angle of $\theta$, is as

5

follows:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos(\theta), & -\sin(\theta) \\ \sin(\theta), & \cos(\theta) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \qquad (2)$$

Where $\cos(\theta)$ and $\sin(\theta)$ are respectively the cosine and sine values of the rotation angle $\theta$. The matrix multiplication operation rotates a two-dimensional point at coordinates $(x, y)$ counterclockwise around the origin point by an angle of $\theta$, yielding the rotated point $(x', y')$.

## 3.2 SLR

### 3.2.1 Keypoint decoupling

We noted that the various components of the keypoint sequences within the same sign language sequence should convey the same semantic information. Thus, we divide the keypoint sequences into four sub-sequences: left hand, right hand, facial expressions and overall, and process them independently. Markers of different colors represent distinct keypoint sequences, as illustrated in Fig 2. This segmentation helps the model more accurately capture the relationships between different parts, facilitating the provision of richer diversity of information. By handling them separately, the model can more attentively capture their respective key features. This keypoint decoupling strategy result and enhances over SLR predictions as shown in our experiments.

### 3.2.2 Keypoint attention module

We employ HRNet Wang et al (2020), which has been trained on COCO-WholeBody Jin et al (2020) dataset, to generate 133 keypoints, including hand, mouth, eye, and body trunk keypoints. Consistent with Chen et al (2022b), we employ a subset of 79 keypoints, comprising 42 hand keypoints, 11 upper body keypoints covering shoulders, elbows, and wrists, and a subset of facial keypoints (10 mouth keypoints and 16 others). Concretely, denoting the keypoint sequence as a multidimensional array with dimensions $C \times T \times N$, where the elements of $C$ consist of $[x_t^n, y_t^n, c_t^n]$, $(x_t^n, y_t^n)$ and $c_t^n$ denotes the coordinates and confidence of the $n$-th keypoint in the $t$-th frame, $T$ denotes the frame number, and $N$ is the total number of keypoints.

As the attention modules for each stream are analogous, we choose the body keypoint attention module as an example for detailed elucidation. The complete attention module is depicted in the Fig. 3. The procedure within the green rounded rectangle outlines the process of single-head attention computation. The input $\boldsymbol{X} \in \mathbb{R}^{C \times T \times N}$ is first enriched with spatial positional encodings. It is then embedded into two linear mapping functions to obtain $\boldsymbol{X} \in \mathbb{R}^{C_e \times T \times N}$, where $C_e$ is usually smaller than $C_{out}$ to alleviate feature redundancy and reduce computational complexity. The attention map is subjected to spatial global normalization. Note that when computing the attention map, we use the Tanh activation function instead of the softmax used in Vaswani et al (2017). This is because the output of Tanh is not restricted to positive values, thus allowing for negative correlations and providing more flexibility Shi et al (2020). Finally, the attention map is element-wise multiplied with the original input to obtain the output features.

To facilitate the model to jointly attend to information from different representation subspaces, the module performs attention computation with $h$ heads. The outputs of all heads are concatenated and mapped to the output space. Similar to the Vaswani et al (2017), we add a feedforward layer at the end to generate the final output. We choose to use leaky ReLU Maas et al (2013) as the non-linear activation function. Additionally, the module includes two residual connections to stabilize network training and integrate different features, as illustrated in the Fig. 3. Finally, we employ 2D convolution to extract temporal features. All processes within the blue rounded rectangle constitute a complete keypoint attention module. It is worth noting that the weights of different keypoint attention modules are not shared.

### 3.2.3 Position encoding

The keypoint sequences are structured into a tensor and inputted to the neural network. Because there is no predetermined sequence or structure for each element of the tensor, we require a positional encoding mechanism to provide a unique label for every joint. Following Vaswani et al (2017); Shi et al (2020), we employ sinusoidal and cosine functions with different frequencies as

encoding functions:

$$PE(p, 2i) = \sin(p/10000^{2i/C_{in}})$$
$$PE(p, 2i+1) = \cos(p/10000^{2i/C_{in}}) \qquad (3)$$

Where $p$ represents the position of the element and $i$ denotes the dimension of the positional encoding vector. Incorporating positional encoding allows the model to capture positional information of elements in the sequence. Their periodic nature provides different representations for distinct positions, enabling the model to better understand the relative positional relationships between elements in the sequence. Joints within a single frame are sequentially encoded, while identical joint across various frames shares a common encoding. It's worth noting that in contrast to the approach proposed in Shi et al (2020), we only introduce positional encoding for the spatial dimension. We use 2D convolution to extract temporal features, eliminating the need for additional temporal encoding as the continuity of time is already considered in the convolution operation.

### 3.2.4 Spatial global regularization

For action detection tasks on skeleton data, the fundamental concept is to utilize known information, namely that each joint of the human body have unique physical or semantic attributes that remain invariant and consistent across all time frames and instances of data. Utilizing this known information, the objective of spatial global regularization is to encourage the model to grasp broader attention patterns, thus better adapting to diverse data samples. This method is achieved by implementing a global attention matrix, presented in the form of $N \times N$, showing the universal relationships among the body joints. This global attention matrix is shared across all data instances and optimized alongside other parameters during training of the network.

### 3.2.5 Head Network

The output feature from the final attention block undergoes spatial pooling to reduce its dimensions to $T/4 \times 256$ before being inputted into the head network in the Fig. 2. The primary objective of the head network is to further capture temporal context. It is comprised of a temporal linear layer, a batch normalization layer, a ReLU layer, along with a temporal convolutional block containing two temporal convolutional layers with a kernel size of 3 and a stride of 1, followed by a linear translation layer and another ReLU layer. The resulting feature, known as gloss representation, has dimensions of $T/4 \times 512$. Subsequently, a linear classifier and a softmax function are utilized to extract gloss probabilities. We use connectionist temporal classification (CTC) loss $\mathcal{L}_{CTC}^{body}$ to optimize the body attention module.

### 3.2.6 Fuse Head and Ensemble

Every keypoint attention module possesses a distinct array of network heads. To thoroughly harness the capabilities of our multi-stream architecture, we integrate an auxiliary fuse head, designed to assimilate outputs from various streams. This fusion head's configuration mirrors that of its counterparts, like the body head, and is likewise governed by CTC loss, represented as $\mathcal{L}_{CTC}^{fuse}$. The forecasted frame gloss probabilities are averaged and subsequently furnished to an ensemble to fabricate the gloss sequence. This ensemble approach amalgamates outcomes from multiple streams, thereby enhancing predictions, as demonstrated in the experiments.

### 3.2.7 Self-Distillation

Frame-Level Self-Distillation Chen et al (2022b) is employed, where the predicted frame gloss probabilities are used as pseudo-targets. In addition to coarse-grained CTC loss, extra fine-grained supervision is provided. Pursuant to our multi-stream design, we use the average gloss probability from the four head networks as pseudo-targets to guide the learning process of each stream. In a formal capacity, we endeavor to diminish the KL divergence between the pseudo-targets and the predictions of the four head networks. This process is designated as frame-level self-distillation loss, for it provides not merely frame-specific oversight but also filters insights from the concluding ensemble into each distinct stream.

### 3.2.8 Loss Function

The overall loss of MSKA-SLR is composed of two parts:1) the CTC losses applied on the outputs of the left stream($\mathcal{L}_{CTC}^{left}$), right stream($\mathcal{L}_{CTC}^{right}$), body

stream($\mathcal{L}_{CTC}^{body}$), fuse stream($\mathcal{L}_{CTC}^{fuse}$); 2) the distillation loss ($\mathcal{L}_{Dist}$). We formulate the recognition loss as follows:

$$L_{SLR} = L_{CTC}^{left} + L_{CTC}^{right} + L_{CTC}^{body} + L_{CTC}^{fuse} + L_{Dist} \tag{4}$$

Up to now, we have introduced all components of MSKA-SLR. Once the training is finished, MSKA-SLR is capable of predicting a gloss sequence by fuse head network.

## 3.3 SLT

The traditional methodologies from previous times frequently described sign language translation (SLT) tasks as challenges in neural machine translation (NMT), where the input to the translation network is visual information. This research followed to this approach and implemented a multi-layer perceptron (MLP) with two hidden layers into the MSKA-SLR framework proposed, followed by the translation process, thereby accomplishing SLT. The network constructed in this manner is named MSKA-SLT, with its architecture illustrated in Fig. 1(b). We chose to utilize employ mBART Liu et al (2020) as the translation network due to its outstanding performance in cross-lingual translation tasks. To fully exploit the multi-stream architecture we designed, we appended an MLP and a translation network to the fuse head. The input to the MLP consists of encoded features by the fuse head network, namely the gloss representations. The translation loss is a standard sequence-to-sequence cross-entropy loss Vaswani et al (2017). MSKA-SLT includes the recognition loss Eq. 4 and the translation loss represented by $L_T$, as specified in the formula:

$$L_{SLT} = L_{SLR} + L_T \tag{5}$$

## 4 Experiments

### *Implementation Details*

To demonstrate the generalization of our methods, unless otherwise specified, we maintain the same configuration for all experiments. The network employs four streams, with each stream consisting 8 attention blocks, and each block containing 6 attention heads. The output channels are set as follows: $64, 64, 128, 128, 256, 256, 256$ and $256$ respectively. For SLR tasks, we utilize a cosine

annealing schedule over 100 epochs and an Adam optimizer with weight decay set to $1e-3$, and an initial learning rate of $1e-3$. The batch size is set to 8. Following Chen et al (2022a,b), we initialize our translation network with mBART-large-cc25[1] pretrained on CC25[2]. We use a beam width of 5 for both the CTC decoder and the SLT decoder during inference. We train for 40 epochs with an initial learning rate of $1e-3$ for the MLP and $1e-5$ for MSKA-SLR and the translation network in MSKA-SLT. Other hyper-parameters remain consistent with MSKA-SLR. We train our models on one Nvidia 3090 GPU.

### 4.1 Datasets and Evaluation Metrics

#### 4.1.1 Phoenix-2014

Phoenix-2014 Koller et al (2015) is from weather forecast broadcasts aired on the German public TV station PHOENIX over a span of three years. This is a German SLR dataset with a vocabulary size of 1081 for glosses. The dataset comprises 5672, 540, and 629 instances in the training, development and testing set.

#### 4.1.2 Phoenix-2014T

Phoenix-2014T Camgoz et al (2018)is an extension of Phoenix-2014, has ascended as the foremost benchmark for SLR and SLT research in recent years Camgoz et al (2018); Simonyan and Zisserman (2014); Tang et al (2021); Tran et al (2015). It encompasses an array of RGB sign language videos performed by a cadre of nine adept signers using German Sign Language (DGS). These videos are meticulously annotated with sentence-level glosses and accompanied by precise German translations transcribed from spoken news content. The dataset is methodically divided into training, development, and testing subsets, the dataset comprises 7096, 519, and 642 video segments, respectively. With a vocabulary size of 1066 for sign glosses and 2887 for German text, Phoenix-2014T provides a rich resource for SLT research. With all ablation studies conducted using this comprehensive dataset.

---

[1]https://huggingface.co/facebook/mbart-large-cc25
[2]https://commoncrawl.org/

**Table 1**: Comparison with previous works on Sign Language Recognition (SLR). WER is adopted as the evaluation metric. Pre: pre-trained.

| Method | Pre | Phoenix-2014 Dev | Phoenix-2014 Test | Phoenix-2014T Dev | Phoenix-2014T Test | CSL-Daily Dev | CSL-Daily Test |
|---|---|---|---|---|---|---|---|
| **RGB-based** | | | | | | | |
| SubUNets Cihan Camgoz et al (2017) | ✓ | 40.8 | 40.7 | - | - | 41.4 | 41.0 |
| LS-HAN Huang et al (2018) | ✓ | - | - | - | - | 39.0 | 39.4 |
| Hybrid CNN-HMM Koller et al (2018) | ✓ | 31.6 | 32.5 | - | - | - | - |
| DNF Cui et al (2019) | ✓ | 23.8 | 24.4 | - | - | 32.8 | 32.4 |
| CNN-LSTM-HMM Koller et al (2020) | ✗ | 26.0 | 26.0 | 22.1 | 24.1 | - | - |
| FCN Cheng et al (2020) | ✗ | 23.7 | 23.9 | 23.3 | 25.1 | 33.2 | 33.5 |
| Joint-SLRT Camgoz et al (2020b) | ✗ | - | - | 24.6 | 24.5 | 33.1 | 32.0 |
| PiSLTRc-R Xie et al (2021) | ✓ | 23.4 | 23.2 | - | - | - | - |
| SignBT Zhou et al (2021a) | ✓ | - | - | 22.7 | 23.9 | 33.2 | 33.2 |
| VAC Min et al (2021) | ✓ | 21.2 | 22.3 | - | - | - | - |
| STMC Zhou et al (2021b) | ✓ | 21.7 | 20.7 | 19.6 | 21.0 | - | - |
| MMTLB Chen et al (2022a) | ✓ | - | - | 21.9 | 22.5 | - | - |
| C$^2$SLR Zuo and Mak (2022) | ✓ | 20.5 | 20.4 | 20.2 | 20.4 | - | - |
| CorrNet Hu et al (2023b) | ✓ | 18.8 | 19.4 | 18.9 | 20.5 | 30.6 | 30.1 |
| TwoStream-SLR Chen et al (2022b) | ✓ | **18.4** | **18.8** | **17.7** | **19.3** | **25.4** | **25.3** |
| SignBERT+ (+ R) Hu et al (2023a) | ✓ | 19.9 | 20.0 | 18.8 | 19.9 | - | - |
| CTCA Guo et al (2023) | ✓ | 19.5 | 20.1 | 19.3 | 20.3 | 31.3 | 29.4 |
| AdaBrowse+ Hu et al (2023c) | ✓ | 19.6 | 20.7 | 19.5 | 20.6 | 31.2 | 30.7 |
| **Keypoint-based** | | | | | | | |
| TwoStream-SLR Chen et al (2022b) | ✓ | 28.6 | 28.0 | 27.1 | 27.2 | 34.6 | 34.1 |
| SignBERT+ Hu et al (2023a) | ✓ | 34.0 | 34.1 | 32.9 | 33.6 | - | - |
| Ours | ✗ | **21.7** | **22.1** | **20.1** | **20.5** | **28.2** | **27.8** |

### 4.1.3 CSL-Daily

CSL-Daily Zhou et al (2021a) is a recently released dataset for the translation of Chinese Sign Language (CSL), recorded in a studio environment. It encompasses 20654 triplets of (video, gloss, text) enacted by ten unique signers. The dataset delves into diverse subjects such as familial existence, healthcare, and academic milieu. CSL-Daily is composed of 18401, 1077, and 1176 partitions in the training, development and testing sections, correspondingly. The vocabulary size is 2000 for sign glosses and 2343 for Chinese text.

### 4.1.4 Evaluation Metrics

Following previous works Chen et al (2022a); Zhou et al (2021b); Camgoz et al (2020b, 2018); Chen et al (2022b); Hu et al (2023a), we adopt word error rate (WER) for SLR evaluation, and BLEU Papineni et al (2002) and ROUGE-L Lin (2004) to evaluate SLT. Lower WER indicates better recognition performance. For BLEU and ROUGE-L, the higher, the better.

## 4.2 Comparison with State-of-the-art Methods

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

For SLR, we compare our recognition network with state-of-the-art methods on Phoenix-2014, Phoenix-2014T and CSL-Daily, as shown in Table 1. The MSKA-SLR achieves 22.1%, 20.5% and 27.8% WER on the test sets of these three datasets, respectively. Typically, keypoint-based approaches are significantly falling short of RGB-based methods; however, our MSKA-SLR has substantially reduced this disparity.

**Table 2**: Performance comparison of MSKA-SLT with methods for SLT on Phoenix-2014T and CSL-Daily.

| | Phoenix-2014T | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Methods** | **Dev** | | | | | **Test** | | | | |
| | ROUGE | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | ROUGE | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
| **RGB-based** | | | | | | | | | | |
| Sign2Text Cihan Camgoz et al (2018) | 31.80 | 31.87 | 19.11 | 13.16 | 9.94 | 31.80 | 32.24 | 19.03 | 12.83 | 9.58 |
| TSPNet Li et al (2020) | - | - | - | - | - | 34.96 | 36.10 | 23.12 | 16.88 | 13.41 |
| MCT Camgoz et al (2020a) | 45.90 | - | - | - | 19.51 | 43.57 | - | - | - | 18.51 |
| SL-Trans Camgoz et al (2020b) | - | 47.26 | 34.40 | 27.05 | 22.38 | - | 46.61 | 33.73 | 26.19 | 21.32 |
| BN-TIN-Trans Zhou et al (2021a) | 46.87 | 46.90 | 33.98 | 26.49 | 21.78 | 46.98 | 47.57 | 34.64 | 26.78 | 21.68 |
| Joint-SLRT Camgoz et al (2020b) | - | 47.73 | 34.82 | 27.11 | 22.11 | - | 48.47 | 35.35 | 27.57 | 22.45 |
| SimulSLT Yin et al (2021) | 36.04 | 36.01 | 22.60 | 16.05 | 12.39 | 35.13 | 35.92 | 22.70 | 16.03 | 12.27 |
| PiSLTRc-T Xie et al (2021) | 47.89 | 46.51 | 33.78 | 26.78 | 21.48 | 48.13 | 46.22 | 33.56 | 26.04 | 21.29 |
| STMC Zhou et al (2021b) | 48.24 | 47.60 | 36.43 | 29.18 | 24.08 | 46.65 | 46.98 | 36.09 | 28.70 | 23.65 |
| SignBT Zhou et al (2021a) | 50.29 | 51.11 | 37.90 | 29.80 | 24.45 | 49.54 | 50.80 | 37.75 | 29.72 | 24.32 |
| MMTLB Chen et al (2022a) | 53.10 | 53.95 | 41.12 | 33.14 | 27.61 | 52.65 | 53.97 | 41.75 | 33.84 | 28.39 |
| TwoStream-SLT Chen et al (2022b) | 54.08 | **54.32** | **41.99** | **34.15** | **28.66** | 53.48 | **54.90** | **42.43** | **34.46** | **28.95** |
| ConSLT Fu et al (2023) | 47.52 | - | - | - | 24.27 | 47.65 | 51.57 | 38.81 | 30.91 | 25.48 |
| SignBERT+(+ R) Hu et al (2023a) | 51.12 | 51.46 | 38.28 | 30.30 | 24.95 | 50.63 | 52.01 | 39.19 | 31.06 | 25.70 |
| XmDA Ye et al (2023) | 52.42 | - | - | - | 25.86 | 49.87 | - | - | - | 25.36 |
| IP-SLT Yao et al (2023) | **54.43** | 54.10 | 41.56 | 33.66 | 28.22 | **53.72** | 54.25 | 41.51 | 33.45 | 27.97 |
| **Keypoint-based** | | | | | | | | | | |
| Skeletor Jiang et al (2021) | 32.66 | 31.97 | 19.53 | 14.01 | 10.91 | 31.80 | 31.86 | 19.11 | 13.49 | 10.35 |
| TwoStream-SLT Chen et al (2022b) | **53.32** | 53.66 | **41.31** | **33.55** | **28.10** | 53.19 | 54.22 | 41.72 | 33.82 | 28.42 |
| SignBERT+ Hu et al (2023a) | 45.53 | 44.45 | 31.88 | 24.59 | 19.86 | 44.89 | 44.35 | 32.09 | 24.92 | 20.41 |
| Ours | 52.67 | **54.09** | 41.29 | 33.24 | 27.63 | **53.54** | **54.79** | **42.42** | **34.49** | **29.03** |

| | CSL-Daily | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Methods** | **Dev** | | | | | **Test** | | | | |
| | ROUGE | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | ROUGE | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
| **RGB-based** | | | | | | | | | | |
| SL-Trans Camgoz et al (2020b) | 37.06 | 37.47 | 24.67 | 16.86 | 11.88 | 36.74 | 37.38 | 24.36 | 16.55 | 11.79 |
| Joint-SLRT Camgoz et al (2020b) | 44.18 | 46.82 | 32.22 | 22.49 | 15.94 | 44.81 | 47.09 | 32.49 | 22.61 | 16.24 |
| SignBT Zhou et al (2021a) | 49.49 | 51.46 | 37.23 | 27.51 | 20.80 | 49.31 | 51.42 | 37.26 | 27.76 | 21.34 |
| MMTLB Chen et al (2022a) | 53.38 | 53.81 | 40.84 | 31.29 | 24.42 | 53.25 | 53.31 | 40.41 | 30.87 | 23.92 |
| TwoStream-SLT Chen et al (2022b) | **55.10** | **55.21** | **42.31** | **32.71** | **25.76** | **55.72** | **55.44** | **42.59** | **32.87** | **25.79** |
| ConSLT Fu et al (2023) | 41.46 | - | - | - | 14.8 | 40.98 | - | - | - | 14.53 |
| XmDA Ye et al (2023) | 49.36 | - | - | - | 21.69 | 49.34 | - | - | - | 21.58 |
| IP-SLT Yao et al (2023) | 44.33 | 45.26 | 31.77 | 22.87 | 16.74 | 44.09 | 44.85 | 31.50 | 22.66 | 16.72 |
| **Keypoint-based** | | | | | | | | | | |
| TwoStream-SLT Chen et al (2022b) | **54.03** | 54.43 | 41.60 | 31.95 | 25.01 | **55.07** | 55.34 | 42.36 | 32.58 | 25.42 |
| Ours | 53.53 | **55.95** | **42.38** | **32.37** | **25.16** | 54.04 | **56.37** | **42.80** | **32.78** | **25.52** |

Among keypoint-based methods, our method significantly surpasses the most challenging competitor TwoStream-SLR Chen et al (2022b) with 5.9%, 6.7% and 6.3% WER improvement on the testing sets of these three datasets, respectively. Note TwoStream-SLR Chen et al (2022b) and SignBERT+ Hu et al (2023a) utilize pre-trained model that leverage more model parameters and additional resources than MSKA-SLR.

For SLT, we compare our MSKA-SLT with state-of-the-art methods on Phoenix-2014T and CSL-Daily as shown in Tab. 2. We achieved BLEU-4 scores of 29.03 and 25.43 on the test sets of these two datasets, respectively, marking an improvement of 0.61 and 0.1 BLEU-4 scores compared to the keypoint-based methods. Furthermore, our approach on the Phoenix-2014T dataset demonstrated a 0.08 improvement in BLEU-4 score compared to the previous state-of-the-art (SOTA) method.

The results indicate that our MSKA demonstrates significant performance enhancements on SLR and SLT. This highlights the benefit of initially decoupling keypoint sequences for multi-stream attention, followed by aggregating distinct stream feature representations, thereby distinguishing our MSKA from previous SLR and SLT systems.

## 4.3 Ablation Studies

### 4.3.1 Impact of Keypoint Augment

To explore the significance of different data augmentation techniques in SLR endeavors, we methodically implemented each augmentation approach in the training of our models and assessed their efficacy in the context of the

**Table 3**: Study the effects of each component of MSKA-SLR on the Phoenix-2014T SLR task.

| Body | Left | Face | Right | Fuse Head | Distillation | Dev | Test |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| ✓ | | | | | | 25.30 | 25.50 |
| | ✓ | | | | | 57.37 | 57.60 |
| | ✓ | ✓ | | | | 34.90 | 36.26 |
| | | | ✓ | | | 35.28 | 35.58 |
| | ✓ | ✓ | | | | 25.86 | 26.05 |
| | ✓ | ✓ | ✓ | | | 25.27 | 25.69 |
| | ✓ | ✓ | ✓ | ✓ | | 23.08 | 23.67 |
| ✓ | ✓ | ✓ | ✓ | ✓ | | 22.69 | 22.70 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **20.09** | **20.54** |

Phoenix-2014T SLR challenge. The outcomes are detailed in Table 4a. It is discernible that the efficacy of model began to diminish upon integrating of translation and scaling data augmentation methods. We posit that these particular augmentation strategies introduced discrepancies in the data alignment with the validation set, consequently resulting in overfitting. Hence, we made the decision to exclusively employ temporal and rotational augmentations.

### 4.3.2 Impact of Each Component

We initially demonstrate the impacts of each stream of MSKA-SLR in Table. 3. In the absence of the multi-stream architecture, the solitary body stream (where one keypoint attention module manages all keypoints) achieves 25.30% and 25.50% WER on the Phoenix-2014T. Within Table. 3, we present the results separately for the left, face and right streams, as well as the fused outcome. This signifies that the precision of segregated streams is substandard compared to that of the solitary body stream, attributable to the loss of certain information. Nonetheless, owing to the distinct focuses and mutual enhancement among these three streams, their fusion culminates in a WER performance of 23.67%, marking a 1.83% enhancement over the solitary body stream. To optimize the attributes that our model attends to, we integrate the body stream into the fusion head, resulting in a WER performance of 22.70%. Ultimately, by incorporating the self-distillation, our framework achieves the optimal outcome, yielding a WER of 20.54%.

Moreover, in our experiments, we also found that in sign language, the right hand exhibits a more prevailing role compared to the left hand. In our study, the results from only using the left hand and the right hand differ by approximately 22%. This discrepancy may be attributed to the fact that in the majority of individuals, the right hand is the dominant hand, while the left hand is the non-dominant hand. Consequently, the right hand is more suitable for performing the detailed and sophisticated gestures essential for sign language. This results in the right hand typically bearing more responsibility and encompassing more information in sign language.

### 4.3.3 Impact of Attention Modules

The influence of network depth on model efficacy stands as a pivotal concern within the realm of deep learning. Broadly speaking, increasing the number of network layers may enhance model performance, but it can also lead to overfitting. Consequently, we have deliberated upon the impact of the number of attention modules on model efficacy. We have designated the number of modules as 6, 8, 10 and 12, as delineated in Table 4b. We ascertain that the pinnacle of performance is attained with 8 modules, yielding the superlative outcome of 20.54% WER. Additionally, we have delved into the ramifications of attention heads within the attention module on the network. This facilitates the model to simultaneously assimilate information across diverse representation subspaces. Each head possesses the capability to concentrate on distinct segments of the input sequence, thereby significantly amplifying the model's eloquent capacity and its adeptness in capturing intricate relationships. To scrutinize

**Table 4**: Ablation studies of: (a) methods for augmenting keypoint data; (b) the impact of varying the number of keypoint attention modules; (c) the effects of varying the number of attention heads in an attention module; (d) the weight of the distillation loss; (e) impact of Spatial Global Regularization; (f) effectiveness of the spatial-temporal attention.

| temporal | rotate | translate | scale | Dev | Test |
|---|---|---|---|---|---|
| | | | | 24.98 | 25.75 |
| ✓ | | | | 20.25 | 21.58 |
| ✓ | ✓ | | | **20.09** | **20.54** |
| ✓ | ✓ | ✓ | | 20.37 | 21.01 |
| ✓ | ✓ | ✓ | ✓ | 21.21 | 22.21 |

(a) The combined effects of various data augmentation techniques.

| Modules | Dev | Test |
|---|---|---|
| 6 | 20.79 | 21.28 |
| 8 | **20.09** | **20.54** |
| 10 | 20.87 | 21.53 |
| 12 | 22.45 | 22.71 |

(b) The number of attention modules.

| Heads | Dev | Test |
|---|---|---|
| 2 | 20.97 | 21.42 |
| 4 | 20.47 | 21.39 |
| 6 | **20.09** | **20.54** |
| 8 | 20.20 | 21.53 |

(c) The number of attention heads.

| L | Dev | Test |
|---|---|---|
| 0.1 | 20.60 | 20.96 |
| 0.5 | 20.12 | 21.08 |
| 1 | **20.09** | **20.54** |
| 2 | 20.77 | 21.79 |

(d) The weight of the distillation loss.

| SGR | Dev | Test |
|---|---|---|
| ✗ | 21.24 | 21.15 |
| ✓ | **20.09** | **20.54** |

(e) SGR: Spatial Global Regularization.

| Spatial-attn | Temporal-attn | Dev | Test |
|---|---|---|---|
| ✓ | | **20.09** | **20.54** |
| ✓ | ✓ | 25.45 | 25.73 |

(f) Spatial-attnention and temporal-attention.

the significance of the number of heads in keypoint attention, we employ assorted quantities of heads and evaluate their performance in the SLR task, as delineated in Table 4c.

### 4.3.4 Impact of Self-Distillation weight

As different streams embody the same meaning, we integrate self-distillation loss at the end of the model to integrate the features learned by each component. It is a hyper-parameter that is designed to balance the effect of CTC loss and the self-distillation loss. We conduct experiments by varying the weight. Table 4d shows that our MSKA-SLR attains the best performance when the weight is set to 1.0.

### 4.3.5 Impact of Spatial Global Regularization

SGR operates on the attention maps within the attention module to mitigate overfitting. In our experiments, delineated in Table 4e, we initially attained a performance of 21.15% WER on the SLR task without incorporating SGR. Subsequently, through the inclusion of SGR, we achieved the optimal performance of

20.54% WER. Moreover, we explored methodologies for managing temporal information in keypoint sequences: 1) reorganizing temporal data via temporal attention post spatial attention, and 2) exclusively employing 2D convolutions devoid of temporal attention. The outcomes are delineated in Table 4f. It is evident that model achieves WER of 25.73% with the inclusion of temporal attention, whereas utilizing only 2D convolutions results in 20.54% WER. This could be ascribed to the augmentation in parameter quantity, the comparatively diminutive dataset extent, and the heightened vulnerability of model to overfitting.

## 5 Conclusion

In this paper, we concentrate on how to introduce domain knowledge into sign language understanding. To achieve the goal, we present a innovative framework named MSKA-SLR which adopts four streams to keypoint sequences for sign language recognition. A variety of methodologies to make the four streams interact with each other. We further extend MSKA-SLR to a sign language translation model by attaching an MLP and a translation network, resulting in the translation framework named MSKA-SLT. Our MSKA-SLR and

MSKA-SLT achieve encouraging improved performance on SLR and SLT tasks across a series of datasets including Phoenix-2014, Phoenix-2014T, and CSL-Daily. We achieved state-of-the-art performance in the Phoenix-2014T sign language translation task. We hope that our approach can serve as a baseline to facilitate future research.

**Data Availability.** The Phoenix-2014 and Phoenix-2014T datasets are publicly available at https://www.i6.informatik. rwth-aachen.de/~koller/RWTH-PHOENIX/ and https://www.i6.informatik.rwth-aachen. de/~koller/RWTH-PHOENIX-2014-T/, respectively. The CSL-Daily datasets will be made available on reasonable request at http://home. ustc.edu.cn/~zhouh156/dataset/csl-daily/.

# References

Bungeroth J, Ney H (2004) Statistical sign language translation. In: Workshop on representation and processing of sign languages, The International Conference on Language Resources and Evaluation, pp 105–108

Camgoz NC, Hadfield S, Koller O, et al (2018) Neural sign language translation. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 7784–7793

Camgoz NC, Koller O, Hadfield S, et al (2020a) Multi-channel transformers for multi-articulatory sign language translation. In: Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16, Springer, pp 301–319

Camgoz NC, Koller O, Hadfield S, et al (2020b) Sign language transformers: Joint end-to-end sign language recognition and translation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp 10023–10033

Carreira J, Zisserman A (2017) Quo vadis, action recognition? a new model and the kinetics dataset. In: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 6299–6308

Chen Y, Wei F, Sun X, et al (2022a) A simple multi-modality transfer learning baseline for sign language translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 5120–5130

Chen Y, Zuo R, Wei F, et al (2022b) Two-stream network for sign language recognition and translation. Advances in Neural Information Processing Systems 35:17043–17056

Cheng KL, Yang Z, Chen Q, et al (2020) Fully convolutional networks for continuous sign language recognition. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIV 16, Springer, pp 697–714

Cihan Camgoz N, Hadfield S, Koller O, et al (2017) Subunets: End-to-end hand shape and continuous sign language recognition. In: Proceedings of the IEEE international conference on computer vision, pp 3056–3065

Cihan Camgoz N, Hadfield S, Koller O, et al (2018) Neural sign language translation. pp 7784–7793

Cui R, Liu H, Zhang C (2019) A deep neural framework for continuous sign language recognition by iterative training. IEEE Transactions on Multimedia 21(7):1880–1891

Dai Z, Yang Z, Yang Y, et al (2019) Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint arXiv:190102860

Feichtenhofer C, Fan H, Malik J, et al (2019) Slowfast networks for video recognition. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 6202–6211

Fu B, Ye P, Zhang L, et al (2023) A token-level contrastive framework for sign language translation. In: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 1–5

Guo L, Xue W, Guo Q, et al (2023) Distilling cross-temporal contexts for continuous sign language recognition. In: Proceedings of the IEEE/CVF conference on computer vision and

pattern recognition, pp 10771–10780

Hu H, Zhao W, Zhou W, et al (2023a) Signbert+: Hand-model-aware self-supervised pre-training for sign language understanding. IEEE Transactions on Pattern Analysis and Machine Intelligence

Hu L, Gao L, Liu Z, et al (2023b) Continuous sign language recognition with correlation network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 2529–2539

Hu L, Gao L, Liu Z, et al (2023c) Adabrowse: Adaptive video browser for efficient continuous sign language recognition. In: Proceedings of the 31st ACM International Conference on Multimedia, pp 709–718

Huang J, Zhou W, Zhang Q, et al (2018) Video-based sign language recognition without temporal segmentation. In: Proceedings of the AAAI Conference on Artificial Intelligence

Jiang T, Camgoz NC, Bowden R (2021) Skeletor: Skeletal transformers for robust body-pose estimation. In: CVPR Workshop, pp 3394–3402

Jin S, Xu L, Xu J, et al (2020) Whole-body human pose estimation in the wild. In: European Conference on Computer Vision, Springer, pp 196–214

Koller O, Forster J, Ney H (2015) Continuous sign language recognition: Towards large vocabulary statistical recognition systems handling multiple signers 141:108–125

Koller O, Zargaran S, Ney H (2017) Re-sign: Re-aligned end-to-end sequence modelling with deep recurrent CNN-HMMs. pp 4297–4305

Koller O, Zargaran S, Ney H, et al (2018) Deep sign: Enabling robust statistical continuous sign language recognition via hybrid CNN-HMMs 126(12):1311–1325

Koller O, Camgoz C, Ney H, et al (2020) Weakly supervised learning with multi-stream CNN-LSTM-HMMs to discover sequential parallelism in sign language videos 42(9):2306–2320

Li D, Xu C, Yu X, et al (2020) Tspnet: Hierarchical feature learning via temporal semantic pyramid for sign language translation. Advances in Neural Information Processing Systems 33:12034–12045

Lin CY (2004) Rouge: A package for automatic evaluation of summaries. In: Text summarization branches out, pp 74–81

Liu Y, Gu J, Goyal N, et al (2020) Multilingual denoising pre-training for neural machine translation. Transactions of the Association for Computational Linguistics 8:726–742

Maas AL, Hannun AY, Ng AY, et al (2013) Rectifier nonlinearities improve neural network acoustic models. In: Proc. icml, Atlanta, GA, p 3

Min Y, Hao A, Chai X, et al (2021) Visual alignment constraint for continuous sign language recognition. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp 11542–11551

Niu Z, Mak B (2020) Stochastic fine-grained labeling of multi-state sign glosses for continuous sign language recognition. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16, Springer, pp 172–186

Papadimitriou K, Potamianos G (2020) Multimodal sign language recognition via temporal deformable convolutional sequence learning. In: Interspeech, pp 2752–2756

Papineni K, Roukos S, Ward T, et al (2002) Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pp 311–318

Pu J, Zhou W, Li H (2019) Iterative alignment network for continuous sign language recognition. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp 4165–4174

Shi L, Zhang Y, Cheng J, et al (2020) Decoupled spatial-temporal attention network for skeleton-based action-gesture recognition. In: Proceedings of the Asian conference on computer vision

Simonyan K, Zisserman A (2014) Two-stream convolutional networks for action recognition in videos. Advances in neural information processing systems 27

Starner T, Weaver J, Pentland A (1998) Real-time American sign language recognition using desk and wearable computer based video. IEEE Transactions on pattern analysis and machine intelligence 20(12):1371–1375

Tamura S, Kawasaki S (1988) Recognition of sign language motion images. Pattern recognition 21(4):343–353

Tang S, Guo D, Hong R, et al (2021) Graph-based multimodal sequential embedding for sign language translation. IEEE Transactions on Multimedia pp 1–1. https://doi.org/10.1109/TMM.2021.3117124

Tran D, Bourdev L, Fergus R, et al (2015) Learning spatiotemporal features with 3D convolutional networks. In: Proceedings of the IEEE international conference on computer vision, pp 4489–4497

Vaswani A, Shazeer N, Parmar N, et al (2017) Attention is all you need. Advances in neural information processing systems 30

Wang J, Sun K, Cheng T, et al (2020) Deep high-resolution representation learning for visual recognition. IEEE transactions on pattern analysis and machine intelligence 43(10):3349–3364

Xie P, Zhao M, Hu X (2021) Pisltrc: Position-informed sign language transformer with content-aware convolution. IEEE Transactions on Multimedia 24:3908–3919

Xie S, Sun C, Huang J, et al (2018) Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In: Proceedings of the European conference on computer vision (ECCV), pp 305–321

Yao H, Zhou W, Feng H, et al (2023) Sign language translation with iterative prototype. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp 15592–15601

Ye J, Jiao W, Wang X, et al (2023) Cross-modality data augmentation for end-to-end sign language translation. arXiv preprint arXiv:230511096

Yin A, Zhao Z, Liu J, et al (2021) SimulSLT: End-to-end simultaneous sign language translation. pp 4118–4127

Zhou H, Zhou W, Li H (2019) Dynamic pseudo label decoding for continuous sign language recognition. In: 2019 IEEE International Conference on Multimedia and Expo (ICME), IEEE, pp 1282–1287

Zhou H, Zhou W, Qi W, et al (2021a) Improving sign language translation with monolingual data by sign back-translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 1316–1325

Zhou H, Zhou W, Zhou Y, et al (2021b) Spatial-temporal multi-cue network for sign language recognition and translation. IEEE Transactions on Multimedia 24:768–779

Zuo R, Mak B (2022) C2SLR: Consistency-enhanced continuous sign language recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 5131–5140

Zuo R, Mak B (2024) Improving continuous sign language recognition with consistency constraints and signer removal. ACM Transactions on Multimedia Computing, Communications and Applications