

PiSLTRc: Position-informed Sign Language Transformer with Content-aware Convolution

Pan Xie, Mengyi Zhao, Xiaohui Hu

Abstract—Since the superiority of Transformer in learning long-term dependency, the sign language Transformer model achieves remarkable progress in Sign Language Recognition (SLR) and Translation (SLT). However, there are several issues with the Transformer that prevent it from better sign language understanding. The first issue is that the self-attention mechanism learns sign video representation in a frame-wise manner, neglecting the temporal semantic structure of sign gestures. Secondly, the attention mechanism with absolute position encoding is direction and distance unaware, thus limiting its ability. To address these issues, we propose a new model architecture, namely PiSLTRc, with two distinctive characteristics: (i) content-aware and position-aware convolution layers. Specifically, we explicitly select relevant features using a novel content-aware neighborhood gathering method. Then we aggregate these features with position-informed temporal convolution layers, thus generating robust neighborhood-enhanced sign representation. (ii) injecting the relative position information to the attention mechanism in the encoder, decoder, and even encoder-decoder cross attention. Compared with the vanilla Transformer model, our model performs consistently better on three large-scale sign language benchmarks: PHOENIX-2014, PHOENIX-2014-T and CSL. Furthermore, extensive experiments demonstrate that the proposed method achieves state-of-the-art performance on translation quality with +1.6 BLEU improvements.

Index Terms—sign language recognition, sign language translation, content-aware neighborhood gathering, position-informed convolution, relative position encoding.

I. INTRODUCTION

SIGN language (SL) is a native language of people with disabled hearing. As a visual language, it consists of various hand gestures, movements, facial expressions, transitions, etc. Sign Language Recognition (SLR) and Translation (SLT) aim at converting the video-based sign languages into sign gloss sequences and spoken language sentences, respectively. Most previous works in this field focus on continuous SLR with the gloss supervision [1–13], few attempts have been made for SLT [14–17]. The main difference is that gloss labels are in the same order with sign gestures, and thus the gloss annotations significantly ease the syntactic alignment under the SLR methods. However, the word ordering rules in natural language are distinct from their counterparts in video-based sign languages [18]. Moreover, sign videos are composed of continuous sign gestures represented by sub-video clips without explicit boundaries. Therefore, directly learning the

mapping between frame-wise signs and natural language words is challenging.

To achieve better translation performance, a promising research line is to perform joint sign language recognition and translation model, which recognizing glosses and translating natural language sentences simultaneously [14, 15]. By doing so, learning with the glosses supervision can better understand sign videos and bring significant benefits to sign language translation. Along this line, Camgzet *et al.* [15] proposes a joint model, Sign Language Transformer (SLTR), which is based on vanilla Transformer [19]. They learn recognition and translation simultaneously and achieve state-of-the-art results due to the Transformer’s advantage in sequence modeling tasks. However, there are still some inherent flaws that limit the capabilities of the Transformer model when solving the SLR and SLT tasks:

- (a) The self-attention mechanism aggregates temporal sign visual features in a frame-wise manner. This mechanism neglects the temporal structure of sign gestures represented by sub-videos, leading to substantial ambiguity in recognition and translation.
- (b) The attention mechanism is permutation-insensitive. Thus position encoding is essential to inject position information for sequence learning, e.g., sign video learning and sentence learning. However, the absolute position encoding used in vanilla Sign Language Transformer (SLTR) [15] is demonstrated distance and direction unaware [20, 21], thus limit its ability for better performance.

To remedy this first shortcoming (a), an intuitive idea is to gather neighboring temporal features to enhance the frame-wise sign representation. However, it is difficult to determine the boundaries of a sign gesture and select the surrounding neighbors precisely. In this paper, we propose a Content-aware and Position-aware Temporal Convolution (CPTcn) to learn robust sign representations. We first propose a content-aware neighborhood gathering method to adaptively select the surrounding neighbors. Specifically, we leverage the local consistency of sign gestures. That is to say, adjacent frames that belong to a sign gesture share similar semantics. Accordingly, we dynamically select neighboring features based on the similarities. Then we aggregate the selected features with temporal convolution layers. However, temporal convolution with a limited receptive field is insufficient to capture the position information of the features in the selected region [22]. To alleviate the drawback, we inject position awareness into convolution layers with Relative Position Encoding (RPE). By aggregating with neighboring similar features, our CPTcn mod-

Pan Xie and Mengyi Zhao are with the School of Automation Science and Electrical Engineering, Beihang University, Beijing 100191, China (e-mail: panxie@gmail.com, mengyizhao@buaa.edu.cn).

Xiaohui Hu is with the Science and Technology on Integrated Information System Laboratory, Institute of Software, Chinese Academy of Sciences, Beijing 100191, China (e-mail: hxh@iscas.ac.cn)

ule obtains discriminative sign representations, thus improving the recognition and translation results.

To solve the second issue (b), we inject relative position information into the learning of sign videos and target sentences. Furthermore, we consider the relative position between sign frames and target words. To the best of our knowledge, we are the first trying to model the position relationship between source sequence and target sequence in sequence-to-sequence architectures. There are several existing methods to endow the self-attention mechanism with relative position information [20, 23–26]. In this paper, we adopt the Disentangled Relative Position Encoding (DRPE) [24] in our video-based sign language learning, target sentence learning, and their mapping learning. Note that, different from RPE mentioned above, DRPE contains the correlations between relative position and sign features, which is proven effective to bring improvements [24, 27]. With the distance and direction awareness learning from DRPE, our improved Transformer model learns better feature representations, thus gaining significant improvements.

We call our approach PiSLTRc for "Position-informed Sign Language TRansformer with content-aware convolution". The overview of our model can be seen in Figure 1. The main technical contributions of our work are summarized as follows:

- 1) We propose a content-aware and position-aware CPTcn module to learn neighborhood-enhanced sign features. Specifically, We first introduce a novel neighborhood gathering method based on the semantic similarities. Then we aggregate the selected features with position-informed temporal convolution layers.
- 2) We endow the Transformer model with relative position information. Compared with absolute position encoding, relative position encoding performs better for sign video and natural sentence learning. Furthermore, we are the first to consider the relative position relationship between sign frames and target words.
- 3) Equipped with the proposed two techniques, our model achieves state-of-the-art performance in translation accuracy on the largest public dataset RWTH-PHOENIX-Weather 2014T. Also, we obtain significant improvements in recognition accuracy compared with other RGB-based models on both PHOENIX-2014 and PHOENIX-2014-T dataset.

The remainder of this paper is organized as follows. Section II reviews related works in sign language and position encoding. Section III introduces the architecture of our proposed PiSLTRc model. Section IV provides implementation details on our model, presents a quantitative analysis that provides some intuition as to why our proposed techniques work, and finally presents the experimental results compared with several baseline models.

II. RELATED WORK

A. Sign Language Recognition

Most previous sign language works focus on continuous sign language recognition (cSLR), which is a weakly supervised sequence labeling problem [2]. cSLR aims at transcribing video-based sign language into gloss sequence. With the released of

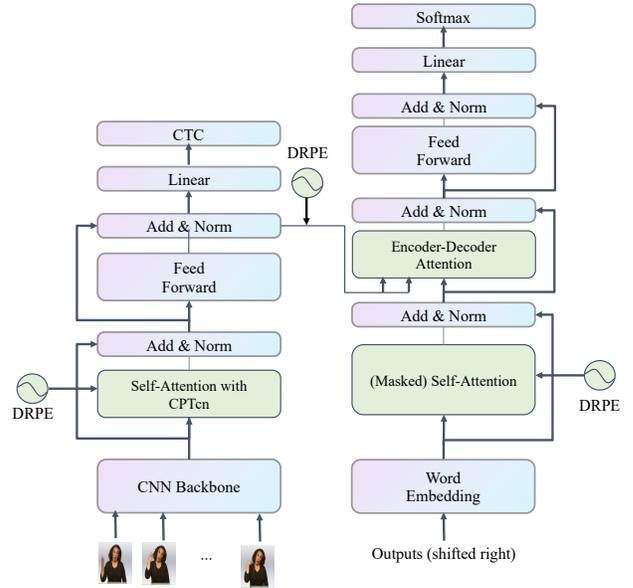


Fig. 1. The overview of our sign language Transformer model equipped with Content-aware and Position-aware Temporal Convolution (CPTcn) and Disentangled Relative Position Encoding (DRPE).

larger-scale cSLR datasets [28], numerous researches burst out implementing sign language recognition tasks in an end-to-end manner [1–13]. The gloss annotations are in same order with sign language, this monotonic relationship significantly ease the syntactic alignment with the cSLR methods. However, the relationship between gloss sequences and the spoken natural language is non-monotonic. Thus it is infeasible to realize SLT with cSLR methods. Fortunately, the knowledge learned by cSLR can be transferred to SLT models and facilitate their performance.

B. Sign Language Translation

Sign language Translation (SLT) is much more challenging because the alignment learning of frame-wise sign gestures and natural language words is difficult. Camgz *et al.* [14] first introduce an end-to-end SLT model that uses Convolution Neural Networks (CNNs) backbone to capture spatial feature and utilizes attention-based encoder-decoder model [29] to learn the mapping of sign videos and natural language sentences. Based on this work, Camgz *et al.* [15] replace the sequence-to-sequence structure with Transformer architecture [19] which is the state-of-the-art model in Neural Machine Translation (NMT) area. Furthermore, they jointly learn the sign language recognition and translation with a shared Transformer encoder and demonstrate that joint training provides significant benefits. Our work is built upon their joint sign language Transformer model, where we improve the Transformer with our proposed CPTcn module and endow the Transformer model with relative position information.

C. Position Encoding in Convolution

Temporal convolution neural network is a common method to model sequential information [30–33]. Convolution layer is

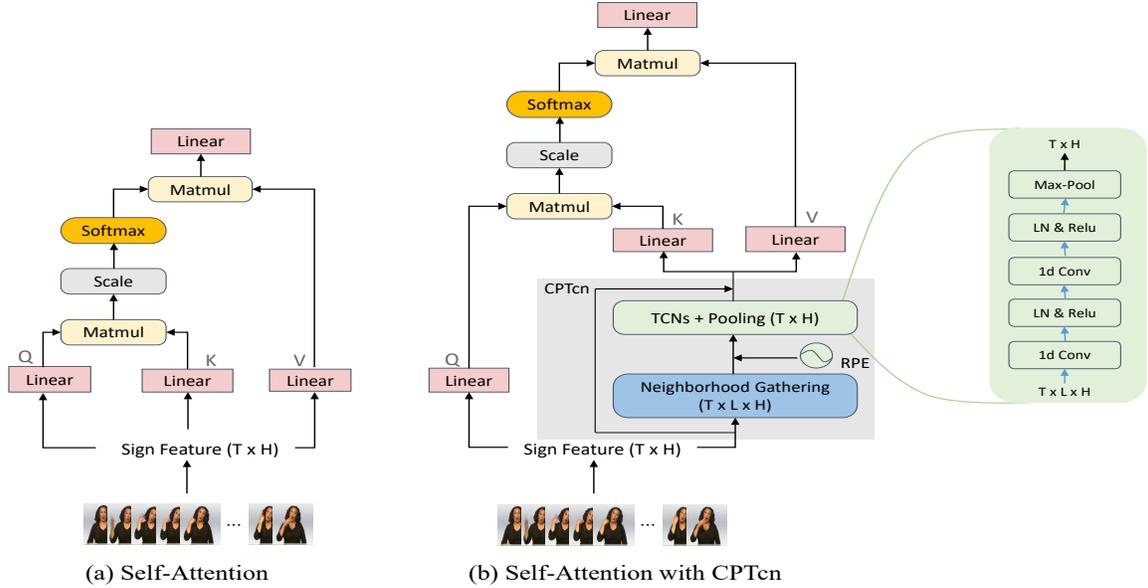


Fig. 2. Vanilla self-attention and self-attention equipped with Content-aware and Position-aware Temporal Convolution (CPTcn). Where neighborhood gathering denotes selecting adjacent relevant features in a contiguous local region, TCNs denote stacked temporal convolution layers.

demonstrated implicitly to learn absolute position information from the commonly used padding operation [22]. However, it is insufficient to learn powerful representations that encode sequential information, especially with the limited receptive field. Explicitly encoding absolute position information is shown effective to learn image features [22]. Upon their hypothesis, we apply relative position encoding (RPE) to the temporal convolution layers, aiming to model the positional correlations between the current feature and its surrounding neighbors.

D. Position Encoding in Self-attention

Transformer entirely relies on the attention mechanism, which does not explicitly model the position information. To remedy the drawback, the sinusoidal absolute position encoding [19] and learnable absolute position encoding [34] are proposed to endow their model with position information. Afterward, relative position encoding is proposed to model long sequence [27] and provides the model with relation awareness [20, 21]. In our work, we reuse the disentangled position encoding [24] to exploit the distance and direction awareness with relative position encoding. Moreover, we also explore the position relationship between sign video and target sentence. Note that, different from RPE in convolution, DRPE in attention mechanism considers the relationship between content and position feature, which is demonstrated effective in previous works [24, 27]. Our experiments indicated that the relative position information is vital for sequence-to-sequence mapping learning.

III. METHOD

A. Preliminaries and Model Overview

Figure 1 illustrates the overall architecture of our proposed model, which jointly learns to recognize and translate sign

videos into gloss annotations and spoken language sentences. In the following subsections, we will first revisit the sign language Transformer structure and then give detailed descriptions about our proposed two methods: content-aware and position-aware temporal convolution (CPTcn), and self-attention with disentangled relative position encoding (DRPE).

B. Joint Sign Language Transformer Structure

Given a series of video frames, the vanilla sign language Transformer (SLTR) model firstly adopts a CNN backbone to extract frame-wise spatial features and uses a word embedding to transfer one-hot natural language words into dense vectors. Then a Transformer-based encoder-decoder model is utilized to learn SLR and SLT simultaneously. For SLR, the encoder output learned temporal sign features. A Connectionist Temporal Classification (CTC) [35] loss is applied to learn the mapping of gloss annotations and sign features. For SLT, the decoder output decomposes sequence level conditional probabilities in an autoregressive manner and then calculates the cross-entropy loss for each word. Meanwhile, the learning of SLR and SLT share the Transformer encoder.

Vanilla Transformer is a sequence-to-sequence structure, which consists of several Transformer blocks. Each block contains a multi-head self-attention and a fully feed-forward network. Given a feature sequence $F \in R^{M \times d}$ with M frames, taking single-head attention as an example, the standard self-attention can be formulated as:

$$\begin{aligned}
 Q &= FW_q, K = FW_k, V = FW_v, \\
 S &= QK^T, \\
 \text{Attn}(Q, K, V) &= \text{softmax}\left(\frac{S}{\sqrt{d}}\right)V, \\
 a_{ij} &= \frac{\exp(s_{ij}/\sqrt{d})}{\sum_{j'} \exp(s_{ij'}/\sqrt{d})}
 \end{aligned} \tag{1}$$

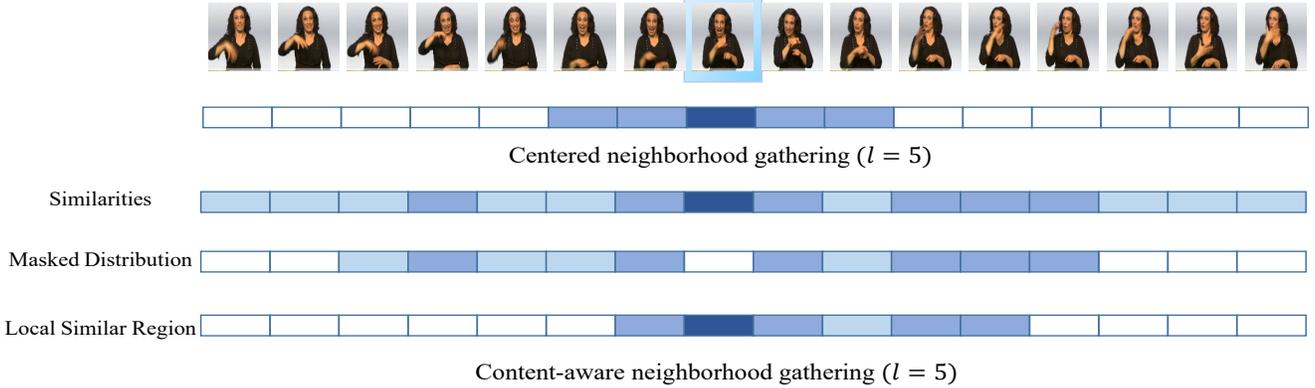


Fig. 3. Our proposed content-aware neighborhood gathering method, compared with the centered neighborhood gathering method.

where $W_q, W_k, W_v \in R^{d \times d}$ represents projection matrices. $S \in R^{M \times M}$ represents the similarity computed by query $Q \in R^{M \times d}$ and key $K \in R^{M \times d}$. a_{ij} represent the normalized attention weights respectively.

Our work concentrates on improving the self-attention mechanism to understand sign video and target sentences better. To focus on our main contributions, we omit the detailed architecture and refer readers to [19] for reference.

C. Content-aware and Position-aware Temporal Convolution

As shown in Figure 2, we propose content-aware and position-aware temporal convolution (CPTcn) to learn local temporal semantics, aiming at obtaining more discriminative sign representations. In this section, we first introduce a content-aware neighborhood gathering method, which adaptively selects surrounding neighbors. Secondly, we elaborate on the detail of endowing the temporal convolution with relative position information, which models the relationship between surrounding features and the current feature. Finally, we incorporate the proposed CPTcn module with the self-attention mechanism.

1) *Content-aware neighborhood gathering Method*: In sign videos, we observe that each sign gesture usually lasts about 0.5~0.6 seconds (~ 16 frames). However, the vanilla Sign Language Transformer (SLTR) model aggregates sign features in a frame-wise manner, thus neglecting the local temporal structure of sign gestures. Unlike their work, we develop a content-aware neighborhood gathering method to adaptively select the relevant surrounding features, which are around a specific feature and in a contiguous region. Shown as Figure 3, we obtain the clip-level feature with neighboring features via three steps:

1). Given the sequential representations $F = \{f_1, f_2, \dots, f_M\}$ from the CNN backbone model, we apply outer tensor product to get a similarity matrix $s \in R^{M \times M}$:

$$s = \frac{FF^T}{\sqrt{d}} \quad (2)$$

where the diagonal elements in s represent similarities towards the features themselves.

2). To ensure neighbors are going to be selected instead of the far-away ones, we only consider a range $[t-l, t+l]$ for a specific feature $f_t \in R^d$ to keep local semantic consistency. Then we replace the similarity scores with *-inf* outside this range and at the current feature. Mathematically, the selecting criterion for f_t becomes as:

$$s_{t,j} = \begin{cases} s_{t,j}, & j \in [\max(0, t-l), t) \ \& \ (t, \min(t+l, M)] \\ -inf, & \text{others} \end{cases} \quad (3)$$

where l represents the maximum distance among the considering features from the current feature. Then we apply the softmax function to obtain the masked distribution $d^t \in R^M$ in the local region around the current feature f_t :

$$d_t = \text{softmax}(s_t) \quad (4)$$

Note that the weight at the current feature is zero, thus the summation of the weights before and after the current feature is 1.

3). It is hard to determine the size and boundaries of the local region. Fortunately, the normalized distribution of similarities obtained in Equation 4 indicates the location of similar neighbors. Therefore, we use the weights of the normalized distribution before and after the current feature to adaptively determine the size of the selected region. Respectively, we define the size before and after the current feature with l_- and l_+ :

$$\begin{aligned} l_- &= \gamma \sum d_{t,j} \cdot l, & \max(0, t-l) \leq j < t \\ l_+ &= \gamma \sum d_{t,j} \cdot l, & t < j \leq \min(t+l, M) \\ l_r &= l_- + l_+ = \gamma \cdot l \end{aligned} \quad (5)$$

where γ is a hyperparameter to control the size of selected region, and the size of the region is l_r . We define the final selected contiguous region as LSR_t (Locally Similar Region) for a specific feature f_t :

$$LSR_t = \{f_{t-l_-}, \dots, f_t, \dots, f_{t+l_+}\} \quad (6)$$

Finally, we adaptively obtain the clip-level features which are in a contiguous local region:

$$f_t^r = \mathcal{F}_{CNG}(f_t, LSR_t) \quad (7)$$

where \mathcal{F}_{CNG} denotes the content-aware neighborhood gathering method, and $f_t^r \in R^{l_r \times d}$ denotes the current feature f_t with its l_r surrounding neighbors. The clip-level features with temporal surrounding neighbors can be computed using Algorithm 1.

Algorithm 1 Content-aware neighborhood gathering Method.

Input: Frame-wise spatial feature $F \in R^{M \times d}$ from CNN backbone;

Output: Clip-level features $F^r \in R^{M \times l_r \times d}$;

```

1:  $s = \frac{1}{\sqrt{d}} f f^T$ ;
2:  $s = \text{local\_mask}(s, l - 1, -\text{inf})$ ;
3:  $s = \text{diagonal\_mask}(s, -\text{inf})$ ;
4:  $a = \text{softmax}(s, \text{dim} = -1)$ ;
5:  $F^r = []$ 
6: for  $t = 0, \dots, M - 1$  do
7:    $l_- = \text{sum}(a[\text{max}(0, t - l) : t]) * l * \gamma$ 
8:    $l_+ = \text{sum}(a[t, t + 1 : \text{min}(t + l, M)]) * l * \gamma$ 
9:    $\text{inds} = [t - l_-, \dots, t, \dots, t + l_+]$ 
10:   $\text{neighbors} = \text{index\_select}(f_t, \text{inds})$ 
11:   $F^r$  append neighbors
12: end for
13: concatenate( $F^r$ , dim = 0)
14: return  $F^r$ ;
```

2) *Position-aware Temporal Convolution*: Temporal convolution is a common method to aggregate sequential features. However, convolution layers with a limited receptive field are insufficient to capture the position information [22], which is important for sign gesture understanding. More specifically, the recognition of sign language is sensitive to the frame order. Absolute position encoding used in previous methods [15, 16] is a promising approach to encode position information. However, it is demonstrated direction- and distance-unaware [21]. Inspired by recent work on language modelling [25], we infuse relative position information to the clip-level feature. We first compute the relative position matrix $R \in R^{M \times l_r}$ between the frame-wise feature and the current feature:

$$R_{t,p} = p - t, \quad t \in [0, M - 1], p \in (t - l_-, t + l_+). \quad (8)$$

Then we represent the relative position indices in learnable embedding space, and obtain the position embeddings $\Phi_{rpe} \in R^{M \times l_r \times d}$. Adding Φ_{rpe} to clip-level features F^r , resulting in position-informed clip-level representation:

$$\hat{F}^r = F^r + \Phi_{rpe}. \quad (9)$$

Lastly, we aggregate the clip-level features with position information to compressed features $F_{ag} \in R^{M \times d}$, and apply a residual function:

$$F_{ag} = \text{MaxPool}(\text{Relu}(\text{LN}(\text{Conv1d}(\hat{F}^r)))) + F \quad (10)$$

where LN represents the Layer Normalization [36]. Conv1d(\cdot) performs a 1-D convolution filter on time dimension with Relu(\cdot) activation function. As shown in Figure 2, we use two layers of such a network, which is omitted here for readability.

3) *Self-attention with CPTcn*: Similar to the vanilla Transformer model, we feed the aggregated feature F_{ag} to the self-attention mechanism. Note that, as shown in Figure 2, we only set $K = V = F_{ag}$, and keep Q as original frame-wise representation F . The reason for this design is to maintain the difference within adjacent features in Q . Experimental result demonstrates that this network design performs better than $Q = K = V = F_{ag}$.

D. Self-Attention with DRPE

As shown in Figure 1, we further inject relative position information into the attention mechanism for sign video learning, target sentence learning, and mapping learning between them. Most existing approaches for endowing the attention mechanism relative position information are based on pairwise distance [20]. They have been explored in machine translation [20], music generation [37] and language modelling [23, 24]. Here, we propose a disentangled relative position encoding (DRPE) [24].

Different from RPE used in Section III-C2, DRPE considers the correlations between relative positions and content features, which are proven that improving the performance [24, 27]. Specifically, we separate the content features and relative position encoding to compute attention weights. The first line of projection in Equation 1 is reparameterized as:

$$\begin{aligned} Q_f &= F W_{q,c}, K_f = F W_{k,c}, V_f = F W_{v,c} \\ Q_p &= P W_{q,p}, K_p = P W_{k,p} \end{aligned} \quad (11)$$

where $F \in R^{M \times d}$ represent the content feature. $Q_f, K_f, V_f \in R^{M \times d}$ represent query, key and value content vectors which are obtained with projection matrices $W_{q,c}, W_{k,c}, W_{v,c} \in R^{d \times d}$. $P \in R^{2\mathcal{L} \times d}$ represents created learned relative position embedding, where \mathcal{L} is the max relative distance. $Q_p, K_p \in R^{2\mathcal{L} \times d}$ represent the projected position embedding with projection matrices $W_{q,p}, W_{k,p} \in R^{d \times d}$, respectively.

Following this, we generate the attention weights with the relative position bias. The calculation of pairwise *content-content* is in the same way as standard self-attention, thereby generating the content-based content vector. While the calculation of pairwise *content-position* is different from standard self-attention. We first create a relative position distance matrix $R^{rel} \in R^{M \times M}$, and then generate the position-based content vectors. The 2 ~ 4 lines of computing attention weights in Equation 1 are reparameterized as:

IV. EXPERIMENTS

A. Dataset and Metrics

We evaluate our method on three datasets, including PHOENIX-2014 [28], PHOENIX-2014-T [14] and Chinese Sign Language (CSL)[38].

PHOENIX-2014 is a publicly available German Sign Language dataset, which is the most popular benchmark for continuous SLR. The corpus was recorded from broadcast news about weather. It contains videos of 9 different signers with a vocabulary size of 1295. The split of videos for Train, Dev, and Test is 5672, 540, and 629, respectively.

PHOENIX-2014-T is the benchmark dataset of sign language recognition and translation. It is an extension of the PHOENIX14 dataset [28]. Parallel sign language videos, gloss annotations, and spoken language translations are available in PHOENIX14T, which makes it feasible to learn SLR and SLT tasks jointly. The corpus is curated from a public television broadcast in Germany, where all signers wear dark clothes and perform sign language in front of a clean background. Specifically, the corpus contains 7096 training samples (with 1066 different sign glosses in gloss annotations and 2887 words in German spoken language translations), 519 validation samples, and 642 test samples.

CSL is a Chinese Sign Language dataset, which is also a widely used benchmark for continuous SLR. These videos were recorded in a laboratory environment, using a Microsoft Kinect camera with a resolution of 1280×720 and a frame rate of 30 FPS. In this corpus, there are 100 sentences, and each sentence is signed five times by 50 signers (in total 2,500 videos). As no official split is provided, we split the dataset by ourselves. We give 20,000 and 5,000 samples to the training set and testing set, respectively. When splitting the dataset, we ensure that the sentences in the training and testing sets are the same, but the signers are different.

We evaluate our model on the performance of SLR and SLT as following [14]:

Sign2gloss aims to transcribe sign language videos to sign glosses. It is evaluated using word error rate (WER), which is a widely used metric for cSLR:

$$\text{WER} = \frac{\#\text{substitution} + \#\text{deletion} + \#\text{insertion}}{\#\text{words in reference}} \quad (13)$$

Sign2text aims to directly translate sign language videos to spoken language translation without intermediary representation. It is evaluated using BLEU [39] which is widely used for machine translation.

Sign2(gloss+text) aims to jointly learn continuous SLR and SLT simultaneously. This approach is currently state-of-the-art in the performance of SLT since the training of cSLR brings benefits for sign video understanding, thus improving the performance of translation.

B. Implementation and Evaluation Details

1) *Network Details*: Like Camgz *et al.* [15], we extract frame-wise spatial features with CNN backbone from CNN-LSTM-HMM [2]. Then we apply the improved Transformer network to learn SLR and SLT simultaneously. Its

$$s_{i,j}^{rel} = \{Q_{f,i}, Q_p R_{i-j}^{rel}\} \times \{K_{f,j}, K_p R_{j-i}^{rel}\}^T$$

$$= \underbrace{Q_{f,i} K_{f,j}^T}_{c2c} + \underbrace{Q_{f,i} K_p^T R_{j-i}^{rel T}}_{c2p} + \underbrace{Q_p R_{i-j}^{rel} K_{f,j}^T}_{p2c} +$$

$$\underbrace{Q_p R_{i-j}^{rel} K_p^T R_{j-i}^{rel T}}_{p2p}$$

$$\text{Attn}(Q_f, K_f, V_f, Q_p R_{q-k}^{rel}, K_p R_{k-q}^{rel}) = \text{softmax}\left(\frac{S^{rel}}{\sqrt{4d}}\right) V_f,$$

$$a_{ij} = \frac{\exp(s_{ij}^{rel} / \sqrt{4d})}{\sum_{j'} \exp(s_{ij'}^{rel} / \sqrt{4d})} \quad (12)$$

where $S^{rel} \in R^{M \times M}$ represents the unnormalized attention score matrix and s_{ij}^{rel} represents the score computed by query at position i and key at the position j . $R_{q-k}^{rel} \in R^{M \times M}$ represents the relative distance matrix computed by the positions of query and key. R_{i-j}^{rel} lies in the (i, j) -th of R_{q-k}^{rel} , and represents the relative distance between i -th query and j -th key. $R_{k-q}^{rel} \in R^{M \times M}$ and R_{j-i}^{rel} are computed in similar ways. Note that R_{i-j}^{rel} and R_{j-i}^{rel} are opposite numbers thus providing our model with directional information.

Moreover, in the first line of the above equation, the first item $c2c$ represents *content-to-content* which is the content-based content vectors. The second and third item $c2p$ and $p2c$ represent *content-to-position* and *position-to-content* respectively, which are relative position based content vectors. $p2p$ represents *position-to-position* which is omitted in vanilla DRPE [24]. However, in our experiments, we find that $p2p$ bring improvements to our performance in both recognition and translation. Therefore, we keep this item of *position-to-position*. In Section IV-C3, we analyze the impact of different item in the first line of Equation 12.

Preceding this, in the last two lines, we apply softmax function and scaling factor $\frac{1}{\sqrt{4d}}$ to get normalized scaled attention weights.

Totally, there are two differences between the DRPE method applied in our architecture and DeBERTa [24]. The first is that we consider the position-to-position information, which is omitted in DeBERTa. Experimental results in Table V show the effectiveness of this item. The second difference is that DRPE is used in text-only in DeBERTa for language modeling. However, in our proposed model, as seen in Figure 1, we apply the relative position method in text-only target sentence learning, image-only sign video learning, and even the cross-modal video sequence and target sentence interaction. Experimental results in Table IV show the effectiveness of our improvements. Note that we are the first to consider the relative position relationship between sign frames and target words.

In summary, equipping with the CPTcn module and DRPE in self-attention layers, the heart module in the Transformer model, we finally arrive at our proposed PiSLTRc model.

TABLE I

EVALUATION OF DIFFERENT NEIGHBORHOOD GATHERING METHOD ON PHOENIX-2014-T. "NG" IS THE ABBREVIATION OF NEIGHBORHOOD GATHERING.

NG method	SLR(WER)		SLT(BLEU-4)	
	DEV	TEST	DEV	TEST
w/o NG	24.23	24.92	20.54	20.80
Centered NG ($l_r = 16$)	23.64	24.17	20.73	21.23
Sparse NG ($l_r = 16$)	23.06	23.52	21.83	22.08
Content-aware NG ($l_r = 16$)	22.23	23.01	23.17	23.40

TABLE II

EVALUATION OF THE SIZE OF THE SELECTED LOCAL SIMILAR REGION ON PHOENIX-2014-T.

Size of LSR	SLR(WER)		SLT(BLEU-4)	
	DEV	TEST	DEV	TEST
$l_r = 8$	22.85	23.85	22.30	22.99
$l_r = 12$	22.52	23.62	22.31	23.51
$l_r = 16$	22.23	23.01	23.17	23.40
$l_r = 20$	23.02	23.74	22.85	22.86

setting used in our experiments is based on Camg z *et al.* [15]. Specifically, we use 512 hidden units, 8 heads, 6 layers, and 0.1 dropout rate.

In our proposed CPTcn model, the size of the select contiguous local similar region l_r is set to be 16 (about 0.5-0.6 seconds), which is the average time needed for completing a gloss. We analyze the impact of the size in Section IV-C1.

The setting of two temporal convolution layers is F3-S1-P0-F3-S1-P0, where F, S, P denote the kernel filter size, stride, and padding size, respectively. The analysis of different modules of the position-informed convolution is concluded in Section IV-C2.

In the self-attention and cross-attention mechanism, we apply DRPE to inject relative position information. We set the max relative distance \mathcal{L} to be 32 in our experiments. The analysis of the DRPE is conduct in Section IV-C3.

Besides, we train the SLR and SLT simultaneously. Thus we set λ_R and λ_T as the weight of recognition loss and translation loss.

2) *Training*: We use the Adam optimizer [40] to optimize our model. We adopt the warmup schedule for learning rate that increases the learning rate from 0 to $6.8e-4$ within the first 4000 warmup steps and gradually decay it with respect to the inverse square root of training steps. We train the model on 1 NVIDIA TITAN RTX GPU, and use 5 checkpoints averaging for the final results.

3) *Decoding*: During inference, we adopt CTC beam search decoder with a beam size of 5 for SLR decoding. Meanwhile, we also utilize the beam search with the width of 5 for SLT decoding, and we apply a length penalty [41] with α values ranging from 0 to 2.

C. Ablation Study

1) *Analysis of content-aware neighborhood gathering method*: In our proposed CPTcn module, we introduce a

TABLE III

EVALUATION OF DIFFERENT MODEL IN CPTcn MODULE ON PHOENIX-2014-T. "PE" DENOTES POSITION ENCODING. "APE" DENOTES ABSOLUTE POSITION ENCODING.

module in CPTcn	SLR (WER)		SLT (BLEU-4)	
	DEV	TEST	DEV	TEST
w/o PE	23.44	24.03	21.15	21.47
w/ APE	22.89	23.47	21.89	22.18
w/ RPE	22.23	23.01	23.17	23.40
w/o Redisual	24.19	25.31	20.78	20.36
w/o LN	23.01	23.72	21.99	21.57
CPTcn	22.23	23.01	23.17	23.40

content-aware neighborhood gathering method to select the relevant surrounding neighbors dynamically. Three potential concerns with using this method are: 1) How many improvements does the content-aware method bring? 2) Must be the selected features contiguous in position? 3) What is the appropriate size of the selected region?

In Table I, we compared three methods to verify the first two questions: the essential of whether the selected region is content-aware and contiguous. For notation, **w/o NG** means no neighborhood gathering method. **Centered NG** means directly to select k features centered around the current feature. **Sparse NG** means dynamically selecting k features with the highest similarity, which may be discontinuous in position. **Content-aware NG** means to select k contiguous features adaptively based on similarity using our proposed content-aware segmentation method. We can see those neighborhood gathering methods effectively improve the performance. **Sparse NG** substantially outperforms **Centered NG**. This gap suggests that the content-aware method is critical for feature selecting. Moreover, **Content-aware NG** performs better than other methods. This indicates that our content-aware contiguous feature aggregation is more suitable for capturing sign gesture representation.

In Table II, we explore the appropriate size of the select local similar region (LSR). The performance of our model performs best when the size of LSR is 16. This is consistent with the finding that the 16-frame (about 0.5-0.6 seconds) is the average time needed for completing a gloss. Besides, by gathering the larger width regions (for example, 20 frames), we observed slight performance degradation. This is because 20 frames (about 1 second) usually contain more than one gesture and thus lower the performance.

2) *Analysis of position-aware Temporal Convolution*: In the first two lines in Table III, we study the relative position encoding in the CPTcn module. Experimental results show that position information is crucial for aggregating the sequential features. Furthermore, compared with absolute position encoding (APE), relative position encoding (RPE) bring improvements with +1.22 BLEU scores and -0.46% WER score on the test dataset. The result supports the conjecture of Yan *et al.* [21] that RPE provides direction and distance awareness for sequence modeling compared with APE method.

Moreover, we explore the network design of the temporal



Fig. 4. Qualitative recognition results of our proposed modules from Dev set (D: delete, I: insert, S: substitute).

feature aggregator method in the 3 ~ 4 lines in Table III. Experiments show that residual connection is essential. And Layer normalization is effective for sequential feature modeling.

3) *Analysis of DRPE in self-attention*: We further conduct comparative experiments to analyze the effectiveness of disentangle relative position encoding (DRPE) in the attention mechanism. As shown in Figure 1, we replace absolute position encoding (APE) with DRPE in three places: encoder self-attention, decoder self-attention, and encoder-decoder cross attention. For notation, in Table IV, "Enc-SA" means self-attention in the encoder module. "Dec-SA" means self-attention in the decoder module. "Enc-Dec-CA" means cross attention between encoder and decoder. In the 1 ~ 4 lines of Table IV, we can see that the DRPE method used in the encoder and decoder all brings significant improvements. This further demonstrates that relative position encoding provides the direction and distance awareness for sequence representation learning. In addition, we find that the performance of DRPE used only in the encoder is better than that of DRPE used only in the decoder. This phenomenon suggests that direction and distance information are more critical for sign video learning than sentence representation learning.

As we move to the fourth line in Table IV, the results show that DRPE in encoder-decoder attention also increases the performance. This phenomenon shows that even if the order of the word in the natural language is inconsistent with the sign language gloss, the relative position information still benefits their mapping learning.

Different from DRPE used in DeBERTa [24], we further explore the effectiveness of different items mentioned in Equation 12 in our task. Experimental results are shown

TABLE IV
ANALYSIS OF DRPE IN ENCODER SELF-ATTENTION AND DECODER SELF-ATTENTION. "SA" IS THE ABBREVIATION OF SELF-ATTENTION. "CA" MEANS THE CROSS-ATTENTION OF ENCODER-DECODER.

Method	SLR (WER)		SLT (BLEU-4)	
	DEV	TEST	DEV	TEST
w/ APE	25.36	25.27	20.12	20.39
Enc-SA w/ DRPE	22.89	23.76	22.35	22.47
Dec-SA w/ DRPE	23.29	23.84	21.89	21.27
Enc-SA & Dec-SA w/ DRPE	22.54	22.89	22.78	22.90
Enc-Dec-CA w/ DRPE	23.74	23.93	21.59	21.41
All w/ DRPE	22.23	23.01	23.17	23.40

TABLE V
ANALYSIS OF DIFFERENT ITEM IN DRPE. "c2c" DENOTES *content-to-content*, "c2p" DENOTES *content-to-position*, "p2c" DENOTES *position-to-content*, AND "p2p" DENOTES *position-to-position*.

Item in DRPE	SLR (WER)		SLT (BLEU-4)	
	DEV	TEST	DEV	TEST
c2c only	25.79	25.85	20.03	20.18
+ c2p & p2c	22.57	23.26	22.84	22.79
+ p2p	22.23	23.01	23.17	23.40

in Table V, the correlations between content and position feature bring significant improvement. Moreover, the *position-to-position* item also benefits our model. This result is consistent with the conclusion in Ke *et al.* [26]. Accordingly, we adopt these four items in our disentangled relative position encoding.

TABLE VI
QUALITATIVE RESULTS WITH DIFFERENT METHODS ON SLT TASK.

GT:	in der nacht sinken die temperaturen auf vierzehn bis sieben grad . (at night the temperatures drop to fourteen to seven degrees .)
SLTR:	heute nacht werte zwischen sieben und sieben grad . (tonight values between seven and seven degrees .)
PiSLTRc:	heute nacht kühlt es ab auf vierzehn bis sieben grad . (tonight it's cooling down to fourteen to seven degrees .)
GT:	an der saar heute nacht milde sechzehn an der elbe teilweise nur acht grad . (on the saar tonight a mild sixteen on the elbe sometimes only eight degrees .)
SLTR:	südlich der donau morgen nur zwölf am oberrhein bis zu acht grad . (south of the danube tomorrow only twelve on the upper rhine up to eight degrees .)
PiSLTRc:	am oberrhein heute nacht bis zwölf am niederrhein nur kühle acht grad . (on the upper rhine tonight until twelve on the lower rhine only a cool eight degrees .)
GT:	am tag von schleswig holstein bis nach vorpommern und zunächst auch in brandenburg gebietsweise länger andauernder regen . (In the south, denser clouds sometimes appear, otherwise it is partly clear or only slightly cloudy .)
SLTR:	am mittwoch in schleswig holstein nicht viel regen . (not much rain on wednesday in schleswig holstein .)
PiSLTRc:	am donnerstag erreicht uns dann morgen den ganzen tag über brandenburg bis zum teil dauerregen . (on thursday we will reach us tomorrow the whole day over brandenburg until partly constant rain.)
GT:	im süden gibt es zu beginn der nacht noch wolken die hier und da auch noch ein paar tropfen fallen lassen sonst ist es meist klar oder nur locker bewölkt . (In the south there are still clouds at the beginning of the night that drop a few drops here and there, otherwise it is mostly clear or only slightly cloudy .)
SLTR:	im süden tauchen im süden teilweise dichtere wolken auf sonst ist es verbreitet klar . (in the south there are sometimes denser clouds in the south otherwise it is widely clear .)
PiSLTRc:	im süden tauchen auch mal dichtere wolken auf sonst ist es gebietsweise klar oder nur locker bewölkt . (In the south, denser clouds sometimes appear, otherwise it is partly clear or only slightly cloudy .)

TABLE VII
THE EVALUATION RESULTS ON SIGN2GLOSS TASK ON PHOENIX-2014-T DATASET.

<i>sign2gloss</i>	DEV		TEST	
	del/ins	WER	del/ins	WER
Model				
DNF [6]	5.9/3.0	22.7	6.8/2.9	23.4
CNN+LSTM+HMM [2]	-	24.5	-	26.5
SLTR-R [15]	-	24.9	-	24.6
FCN [10]	6.5/3.2	22.9	5.8/4.7	23.7
STMC (RGB) [12]	-	25.0	-	-
PiSLTRc-R (ours)	4.9/4.2	21.8	5.1/4.4	22.9

4) *Qualitative Analysis on SLR*: In Figure 4, we show two examples with different methods on the SLR task. Equipped with proposed approaches, our PiSLTRc model learns accurate sign gesture recognition and thus achieving significant improvements. Furthermore, we find that the model trained based on CTC loss function tends to predict "peak" on the continuous gestures. And our proposed CPTcn model is adequate to alleviate this situation. As shown in Figure 4, the recognition of adjacent frames in a contiguous region is more precise.

5) *Qualitative Analysis on SLT*: In Table VI, we show several examples with different models on the SLT task. Compared with the vanilla SLTR model [15], our proposed PiSLTRc produces target sentences with higher quality and accuracy.

Comparing the translation results of the first example as illustrated in Table VI, we see that "vierzehn (fourteen)" is mistranslated as "sieben (seven)" in SLTR model. However, it is correctly translated in our PiSLTRc model. As we move to the

TABLE VIII
THE EVALUATION RESULTS ON SIGN2GLOSS TASK ON PHOENIX-2014 DATASET.

<i>sign2gloss</i>	DEV		TEST	
	del/ins	WER	del/ins	WER
Model				
DeepHand [1]	16.3/4.6	47.1	15.2/4.6	45.1
DeepSign [42]	12.6/5.1	38.3	11.1/5.7	38.8
SubUNets [11]	14.6/4.0	40.8	14.3/4.0	40.7
Staged-Opt [3]	13.7/7.3	39.4	12.2/7.5	38.7
Re-Sign [43]	-	27.1	-	26.8
DNF [6]	7.8/3.5	23.8	7.8/3.4	24.4
CNN-LSTM-HMM [2]	-	26.0	-	26.0
FCN [10]	-	23.7	-	23.9
STMC (RGB) [12]	-	25.0	-	-
SBD-RL [44]	9.9/5.6	28.6	8.9/5.1	28.6
SLTR-R(our implementation)	8.9/4.2	24.5	9.0/4.3	24.6
PiSLTRc-R (ours)	8.1/3.4	23.4	7.6/3.3	23.2

second example in this table, we see that "heute nacht (tonight)" is mistranslated as "morgen (tomorrow)" in SLTR model, and it is correctly in our PiSLTRc model. To sum up, specific numbers and named entities are challenging since there is no grammatical context to distinguish one from another. However, in these two examples, we see that our model translates specific numbers and named entities more precisely. This demonstrates that our proposed model has a stronger ability to understand sign videos.

When we move to the third and fourth example in the Table VI, we see that our model generate complete sentence with less under-translation. For example, in the third example, "gebietsweise länger andauernder regen (rain lasting longer in

TABLE IX
THE EVALUATION RESULTS ON SIGN2TEXT TASK ON ON PHOENIX2014T DATASET.

<i>sign2text</i>	DEV					TEST				
	Model	ROUGE	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE	BLEU-1	BLEU-2	BLEU-3
RNN-based[14]	31.80	31.87	19.11	13.16	9.94	31.80	32.24	19.03	12.83	9.58
TSPnet [16]	-	-	-	-	-	34.96	36.10	23.12	16.88	13.41
SLTR-T [15]	-	45.54	32.60	25.30	20.69	-	45.34	32.31	24.83	20.17
Multi-channel [17]	44.59	-	-	-	19.51	43.57	-	-	-	18.51
PiSLTRc-T (ours)	47.89	46.51	33.78	26.78	21.48	48.13	46.22	33.56	26.04	21.29

TABLE X
THE EVALUATION RESULTS ON SIGN2(GLOSS+TEXT) TASK ON ON PHOENIX2014T DATASET.

<i>sign2(gloss+text)</i>	DEV						TEST					
	Model	WER	ROUGE	BLEU-1	BLEU-2	BLEU-3	BLEU-4	WER	ROUGE	BLEU-1	BLEU-2	BLEU-3
RNN-based[14]	-	44.14	42.88	30.30	23.02	18.40	-	43.80	43.29	30.39	22.82	18.13
SLTR [15]	24.61	-	46.56	34.03	26.83	22.12	24.49	-	47.20	34.46	26.75	21.80
Multi-channel [17]	-	-	-	-	-	22.38	-	-	-	-	-	21.32
STMC(RGB-based) [45]	-	44.30	44.06	32.69	25.45	20.74	-	44.70	45.08	33.80	26.44	21.55
PiSLTRc (ours)	22.23	49.87	47.37	35.41	28.09	23.17	23.01	49.72	48.50	35.97	28.37	23.40

TABLE XI
THE EVALUATION RESULTS ON SIGN2GLOSS TASK ON CSL DATASET.

<i>sign2gloss</i>	WER
S2VT [46]	25.0
LS-HAN [38]	17.3
HLSTM-attn [47]	10.2
CTF [48]	11.2
DenseTCN [49]	14.3
SF-Net [50]	3.8
FCN [10]	3.0
SLTR-R(our implementation)	3.7
PiSLTRc-R (ours)	2.8

some areas)” is under-translated in SLTR model, while it is correctly translated as ”bis zum teil dauerregen (partly constant rain)” in our PiSLTRc model.

In summary, our proposed model performs better than the previous SLTR model when facing the specific numbers and name entities, which are challenging to translate since there is no grammatical context to distinguish one from another. Moreover, the sentences produced follow standard grammar. Nevertheless, it may be improved on the translation quality of the long sentences in the future.

6) *Limitation*: We leverage neighboring similar features to enhance sign representation. The selected features are in a fixed-size region. This is not consistent with the characteristics of sign language. That is to say, the number of frames corresponding to different sign gestures is dynamic.

D. Comparison Against Baselines

In this section, we compare several state-of-the-art models to demonstrate the effectiveness of our work. Similar to Camgiz *et al.* [15], we elaborate the comparison between our proposed

model and baseline models in the three tasks: *sign2gloss*, *sign2text*, and *sign2(gloss+text)*.

1) *sign2gloss*: We evaluate this task in three datasets: PHOENIX-2014-T, PHOENIX-2014 and CSL.

In Table VII, we compare our model with several methods for the *sign2gloss* task on PHOENIX-2014-T dataset. DNF [6] adopt iterative optimization approaches to tackle the weakly supervised problem. They first train an end-to-end recognition model for alignment proposal, and then use the alignment proposal to tune the feature extractor. CNN-LSTM-HMM [2] embeds powerful CNN-LSTM models in multi-stream HMMs and combines them with intermediate synchronization constraints among multiple streams. Vanilla SLTR-R [15] uses the backbone pretrained with CNN-LSTM-HMM setup and then employes a two-layered transformer encoder model. FCN [10] is built upon an end-to-end fully convolutional neural network for cSLR. Furthermore, they introduce a Gloss Feature Enhancement (GFE) to enhance the frame-wise representation, where GFE is trained to provide a set of alignment proposals for the frame feature extractor. STMC (RGB) [12] proposes a spatial-temporal multi-cue network to learn the video-based sequence. For a fair comparison, we only selected the RGB-based model of STMC without leveraging the additional information of hand, face, and body pose. PiSLTRc-R is our model which is trained when the weight of translation loss λ_T is set zero. Similar to vanilla SLTR-R, our work extracts feature from the CNN-LSTM-HMM backbone. As shown in this table, our proposed PiSLTRc-R surpasses the vanilla SLTR model by 12% and 7% on Dev and Test datasets, respectively. Furthermore, in the RGB-based models, we achieve state-of-the-art performance on the *sign2gloss* task.

In Table VIII we also evaluate our PiSLTRc-R model on the PHOENIX-2014 dataset. Compared with existed baseline models, our proposed model achieves comparable results. Note that the vanilla SLTR-R does not report the experimental results on the PHOENIX-2014 dataset. We implement it by ourselves.

Compared with SLTR-R, our PiSLTRc-R model gains 4% and 5% improvements on Dev and Test datasets, respectively.

In Table XI we conduct experiments on CSL dataset. We see that our proposed PiSLTRc-R model achieves state-of-the-art performance. Compared with the SLTR-R model, our PiSLTRc-R model gains 24% improvements on the Test datasets (5,000 examples split by ourselves), respectively.

2) *sign2text*: In Table IX, we compare our approach with several *sign2text* methods on PHOENIX-2014-T dataset. The RNN-based model [14] adopt full frame features from Re-sign. TSPnet [16] utilizes I3D [51] to extract the spatial features, and further finetune I3D on two WSLR datasets [52, 53]. Multi-channel [17] allows both the inter and intra contextual relationship between different asynchronous channels to be modelled within the transformer network itself. PiSLTRc-T is our model that training with the weight of recognition loss λ_R being zero. Like in *sign2gloss*, SLTR-T and our PiSLTRc-T model utilize the pretrained feature from CNN-LSTM-HMM. Experimental results show that our proposed model achieves state-of-the-art performance and surpasses the vanilla SLTR-T model by 3.8% and 5.6% BLEU-4 scores.

3) *sign2(gloss+text)*: In Table X, we compare our model on *sign2(gloss+text)* task. In this task, we jointly learn sign language recognition and translation simultaneously. Namely, λ_R and λ_T are set as non-zero. Note that different settings will obtain different results. Weighing up the performance on recognition and translation in our experiments, we set $\lambda_R = \lambda_T = 1.0$. Compared with vanilla SLTR, our model gains significant improvements on both two tasks. Experiments demonstrate that our proposed techniques bring significant improvements for recognition and translation quality based on the sign language Transformer model.

V. CONCLUSION

In this paper, we indicate two drawbacks of the sign language Transformer (SLTR) model for sign language recognition and translation. The first shortcoming is that self-attention aggregates sign visual features in a frame-wise manner, thus neglecting the temporal semantic structure of sign gestures. To overcome this problem, we propose a CPTcn module to generate neighborhood-enhanced sign features by leveraging the temporal semantic consistency of sign gestures. Specifically, we introduce a novel content-aware neighborhood gathering method to select relevant features dynamically. And then, we apply position-informed temporal convolution layers to aggregate them.

The second disadvantage is the absolute position encoding used in the vanilla SLTR model. It is demonstrated unable to capture the direction and distance information, which are critical for sign video understanding and sentence learning. Therefore, we inject relative position information to SLTR model with disentangled relative position encoding (DRPE) method. Extensive experiments on two large-scale sign language datasets demonstrate the effectiveness of our PiSLTRc framework.

REFERENCES

- [1] 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," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3793–3802, 2016.
- [2] O. Koller, N. C. Camgoz, H. Ney, and R. Bowden, "Weakly supervised learning with multi-stream cnn-lstm-hmms to discover sequential parallelism in sign language videos," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, pp. 2306–2320, 2020.
- [3] R. Cui, H. Liu, and C. Zhang, "Recurrent convolutional neural networks for continuous sign language recognition by staged optimization," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1610–1618, 2017.
- [4] J. Pu, W. Zhou, and H. Li, "Dilated convolutional network with iterative optimization for continuous sign language recognition," in *International Joint Conferences on Artificial Intelligence Organization (IJCAI)*, 2018, pp. 885–891.
- [5] H. Wang, X. Chai, and X. Chen, "A novel sign language recognition framework using hierarchical grassmann covariance matrix," *IEEE Transactions on Multimedia*, vol. 21, pp. 2806–2814, 2019.
- [6] R. Cui, H. Liu, and C. Zhang, "A deep neural framework for continuous sign language recognition by iterative training," *IEEE Transactions on Multimedia*, vol. 21, pp. 1880–1891, 2019.
- [7] J. Pu, W. Zhou, and H. Li, "Iterative alignment network for continuous sign language recognition," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4160–4169, 2019.
- [8] J. Pu, W. Zhou, H. Hu, and H. Li, "Boosting continuous sign language recognition via cross modality augmentation," *Proceedings of the 28th ACM International Conference on Multimedia*, 2020.
- [9] H. Hu, W. Zhou, J. Pu, and H. Li, "Global-local enhancement network for nmfs-aware sign language recognition," *ArXiv*, vol. abs/2008.10428, 2020.
- [10] 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)*. Springer, 2020, pp. 697–714.
- [11] N. C. Camgöz, S. Hadfield, O. Koller, and R. Bowden, "Subunets: End-to-end hand shape and continuous sign language recognition," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 3075–3084, 2017.
- [12] H. Zhou, W. Zhou, Y. Zhou, and H. Li, "Spatial-temporal multi-cue network for continuous sign language recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 34, no. 07, 2020, pp. 13 009–13 016.
- [13] 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*, vol. 29, no. 9, pp. 2822–2832, 2019.

- [14] N. C. Camgöz, S. Hadfield, O. Koller, H. Ney, and R. Bowden, “Neural sign language translation,” *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7784–7793, 2018.
- [15] N. C. Camgöz, O. Koller, S. Hadfield, and R. Bowden, “Sign language transformers: Joint end-to-end sign language recognition and translation,” *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10 020–10 030, 2020.
- [16] 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 *Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [17] N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, “Multi-channel transformers for multi-articulatory sign language translation,” in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 301–319.
- [18] R. Pfau, M. Salzmann, and M. Steinbach, “The syntax of sign language agreement: Common ingredients, but unusual recipe,” *Glossa*, vol. 3, pp. 1–46, 2018.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998–6008.
- [20] P. Shaw, J. Uszkoreit, and A. Vaswani, “Self-attention with relative position representations,” in *2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, (NAACL-HLT)*, 2018, pp. 464–468.
- [21] H. Yan, B. Deng, X. Li, and X. Qiu, “Tener: Adapting transformer encoder for named entity recognition,” *ArXiv*, vol. abs/1911.04474, 2019.
- [22] M. A. Islam, S. Jia, and N. D. B. Bruce, “How much position information do convolutional neural networks encode?” in *International Conference on Learning Representations (ICLR)*, 2020.
- [23] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov, “Transformer-xl: Attentive language models beyond a fixed-length context,” in *Proceedings of the 57th Conference of the Association for Computational Linguistics (ACL)*, 2019.
- [24] P. He, X. Liu, J. Gao, and W. Chen, “Deberta: Decoding-enhanced bert with disentangled attention,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [25] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.
- [26] G. Ke, D. He, and T. Liu, “Rethinking positional encoding in language pre-training,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [27] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” in *Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [28] J. Forster, C. Schmidt, O. Koller, M. Bellgardt, and H. Ney, “Extensions of the sign language recognition and translation corpus rwth-phoenix-weather,” in *Proceedings of the Ninth International Conference on Language Resources and Evaluation, (LREC)*, 2014.
- [29] T. Luong, H. Pham, and C. D. Manning, “Effective approaches to attention-based neural machine translation,” in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2015.
- [30] B. Martinez, P. Ma, S. Petridis, and M. Pantic, “Lipreading using temporal convolutional networks,” in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6319–6323.
- [31] Z. Xia, X. Hong, X. Gao, X. Feng, and G. Zhao, “Spatiotemporal recurrent convolutional networks for recognizing spontaneous micro-expressions,” *IEEE Transactions on Multimedia*, vol. 22, no. 3, pp. 626–640, 2020.
- [32] D. Guo, W. Li, and X. Fang, “Fully convolutional network for multiscale temporal action proposals,” *IEEE Transactions on Multimedia*, vol. 20, no. 12, pp. 3428–3438, 2018.
- [33] S. Zhang, S. Zhang, T. Huang, and W. Gao, “Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching,” *IEEE Transactions on Multimedia*, vol. 20, no. 6, pp. 1576–1590, 2018.
- [34] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, (NAACL-HLT)*, 2019.
- [35] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Machine Learning, Proceedings of the Twenty-Third International Conference, (ICML)*, 2006.
- [36] J. Ba, J. Kiros, and G. E. Hinton, “Layer normalization,” *ArXiv*, vol. abs/1607.06450, 2016.
- [37] C.-Z. A. Huang, A. Vaswani, J. Uszkoreit, I. Simon, C. Hawthorne, N. Shazeer, A. M. Dai, M. Hoffman, M. Dinculescu, and D. Eck, “Music transformer: Generating music with long-term structure,” in *International Conference on Learning Representations (ICLR)*, 2019.
- [38] J. Huang, W. Zhou, Q. Zhang, H. Li, and W. Li, “Video-based sign language recognition without temporal segmentation,” in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [39] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2002.
- [40] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR*, vol. abs/1412.6980, 2015.
- [41] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey,

- J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. S. Corrado, M. Hughes, and J. Dean, "Google's neural machine translation system: Bridging the gap between human and machine translation," *ArXiv*, vol. abs/1609.08144, 2016.
- [42] O. Koller, S. Zargaran, H. Ney, and R. Bowden, "Deep sign: Hybrid cnn-hmm for continuous sign language recognition," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2016.
- [43] O. Koller, S. Zargaran, and H. Ney, "Re-sign: Re-aligned end-to-end sequence modelling with deep recurrent cnn-hmms," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3416–3424, 2017.
- [44] C. Wei, J. Zhao, W. gang Zhou, and H. Li, "Semantic boundary detection with reinforcement learning for continuous sign language recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, pp. 1138–1149, 2021.
- [45] H. Zhou, W. Zhou, Y. Zhou, and H. Li, "Spatial-temporal multi-cue network for sign language recognition and translation," *IEEE Transactions on Multimedia*, pp. 1–1, 2021.
- [46] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, and K. Saenko, "Sequence to sequence – video to text," *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 4534–4542, 2015.
- [47] D. Guo, W. Zhou, H. Li, and M. Wang, "Hierarchical lstm for sign language translation," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [48] S. Wang, D. Guo, W. gang Zhou, Z. Zha, and M. Wang, "Connectionist temporal fusion for sign language translation," *Proceedings of the 26th ACM international conference on Multimedia (ACM MM)*, 2018.
- [49] D. Guo, S. Wang, Q. Tian, and M. Wang, "Dense temporal convolution network for sign language translation," in *International Joint Conferences on Artificial Intelligence Organization (IJCAI)*, 2019.
- [50] Z. Yang, Z. Shi, X. Shen, and Y.-W. Tai, "Sf-net: Structured feature network for continuous sign language recognition," *ArXiv*, vol. abs/1908.01341, 2019.
- [51] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4724–4733, 2017.
- [52] D. Li, C. Rodriguez-Opazo, X. Yu, and H. Li, "Word-level deep sign language recognition from video: A new large-scale dataset and methods comparison," *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1448–1458, 2020.
- [53] H. R. V. Joze and O. Koller, "Ms-asl: A large-scale data set and benchmark for understanding american sign language," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2019.