

1 We thank all the reviewers for their constructive comments. Below, we address concerns raised by reviewers.

2 **In-the-wild experiments (R1, R3).** Because of our insufficient domain knowledge on German sign languages, we
3 opt to collect 76 short videos in German sign language from the internet as in [10]. In this way, we ensure our model
4 trained on RPWT fully covers the vocabularies, thus facilitating the validation of translations. Results show that the
5 proposed TSPNet correctly translates 32 phrases despite the domain gap. Although there are a large number of sign
6 videos, annotating them with glosses requires expert knowledge while being time-consuming and laborious. Without
7 such gloss annotations, our method is still able to construct a sign language translation model. Moreover, the data for
8 training our model are easy to collect similar to [10]. These reasons make our method more favorable.

9 **Discussions on Cihan et al., 2020 (R1, R3).** Cihan et al. employ a visual encoder finetuned on the RWTH-PHOENIX
10 dataset (GSL) with gloss annotations. Training/finetuning on gloss annotations will facilitate a model to effectively
11 infer sign boundaries. In other words, Cihan et al.’s method requires gloss annotations of a particular dataset to achieve
12 its best performance. Thus, comparing our results with theirs is unfair since we do not use gloss annotations.

13 **Current limitations (R1).** (i) We notice that low-frequency words, such as city names, are very challenging to translate,
14 we plan to provide our models with an external knowledge database and allow it to search similar gestures for translation.
15 (ii) Facial expressions often reflect the extent of an event, e.g. shower vs. rain storm, which are currently not well
16 interpreted by our models. Our future work will try to incorporate facial landmarks to improve translation accuracy. We
17 will include the above discussions in the revision.

18 **Alternative visual backbone and language decoder (R2, R3).** (i) TSPNet aims at proposing a generic learning
19 framework to translate continuous signing videos to natural language sentences directly. We adopt I3D as visual feature
20 backbone considering its recent success in sign language interpretation [8, 9, 10]. As SlowFast network demonstrates
21 its superior action modeling capacity to I3D, we believe employing a stronger visual backbone, like SlowFast network,
22 would lead to better translation results. Due to the tight timeframe of rebuttal, we cannot report results using the
23 SlowFast network. We will add the results and discuss the impact of alternative visual encoders in the revised version.
24 (ii) As suggested by R3, after replacing a Transformer decoder with an RNN in [2], we did not observe performance
25 changes in BLEU-4 (+0.09) or ROUGE-L (-0.03). This implies that the performance gains mainly come from the
26 proposed hierarchical encoding scheme rather the decoder.

27 **Clarification on backbones (i) (R2, R4)** We replace Conv2d-based features with I3D features into an RNN model
28 and obtain BLEU-4 of 11.27 and ROUGE-L of 32.22, similar to TSPNet-Single. Since our I3D backbone captures
29 the spatial-temporal clues from video segments, it will leverage neighboring sharp frames to infer motion blurred sign
30 gestures. In addition, our hierarchical feature learning effectively exploits windowing segments at different granularities,
31 facilitating to capture fine-grained gestures. (ii) (R3) We clarify [2] uses AlexNet (see Sec. 5 in [2]), while Koller et al.
32 and Cihan et al. use GoogleNet initialized on ILSVRC and finetuned on RWTH with glosses.

33 **User studies (R3).** Due to the diversity of translations, standard metrics may not fully reflect whether translation results
34 convey the correct meanings or not. To further demonstrate our performance gains, we ask two participants to compare
35 the predictions of TSPNet-Joint and [2] on RPWT to ground-truth translations, and choose the semantically more
36 relevant result. The two participants favor 562 (87.54%) and 580 (90.34%) translations of TSPNet-Joint over those of
37 [2] out of 642 testing instances. This experiment further demonstrates our significant performance improvements.

38 **Difference between SLT and generic action localization (R3).** Since signing progresses continuously, there is no
39 clear boundary between consecutive sign gestures. Thus, it is ambiguous whether each sign gesture or an entire sign
40 sequence should be treated as an action. Moreover, natural languages differ from sign languages grammatically. As
41 a result, segmenting continuous signing into isolated signs is very difficult, especially in the absence of glosses. In
42 contrast, as for the action localization tasks, temporal boundaries of an activity are evidently clear. This results in the
43 essential difference between these two tasks. Although our approach is designed for SLT, we agree with R3 that it
44 would be interesting to see our method generalizes to weakly supervised action localization, where only action labels
45 are provided without boundary annotations. Due to the rebuttal timeframe, we have to leave this in future work.

46 **Computation cost (R3).** TSPNet-Joint and TSPNet-Single share the same architecture while TSPNet-Joint has 18%
47 more operations for hierarchical feature learning over multi-granularity features. TSPNet-Joint takes 110 ms on average
48 to process an RPWT video on an NVIDIA V100 GPU. We will include the discussion in the revision.

49 **Segmenting isolated signs (R4).** Please note that we did not claim our method achieves accurate temporal segmentation.
50 Due to the grammatical differences between sign languages and natural languages, it is very difficult to segment sign
51 gestures especially when glosses are unavailable. This motivates us to explore multi-scale windowing segments to learn
52 representative features for sign language translation in an end-to-end fashion. Furthermore, Table 3 indicates that it is
53 non-trivial to integrate multi-scale segment features. In this regard, our hierarchical feature learning method provides an
54 effective solution to representing gestures for sentence translation in the absence of gloss annotations.