

Sign Language Transformers: Joint End-to-end Sign Language Recognition and Translation

Necati Cihan Camgöz[✉], Oscar Koller[✉], Simon Hadfield[✉] and Richard Bowden[✉]

[✉]CVSSP, University of Surrey, Guildford, UK, [✉]Microsoft, Munich, Germany

{n.camgoz, s.hadfield, r.bowden}@surrey.ac.uk, oscar.koller@microsoft.com

Abstract

Prior work on Sign Language Translation has shown that having a mid-level sign gloss representation (effectively recognizing the individual signs) improves the translation performance drastically. In fact, the current state-of-the-art in translation requires gloss level tokenization in order to work. We introduce a novel transformer based architecture that jointly learns Continuous Sign Language Recognition and Translation while being trainable in an end-to-end manner. This is achieved using a Connectionist Temporal Classification (CTC) loss to bind the recognition and translation problems into a single unified architecture. This joint approach does not require any ground-truth timing information, simultaneously solving two co-dependant sequence-to-sequence learning problems and leads to significant performance gains.

We evaluate the recognition and translation performances of our approaches on the challenging RWTH-PHOENIX-Weather-2014T (PHOENIX14T) dataset. We report state-of-the-art sign language recognition and translation results achieved by our Sign Language Transformers. Our translation networks outperform both sign video to spoken language and gloss to spoken language translation models, in some cases more than doubling the performance (9.58 vs. 21.80 BLEU-4 Score). We also share new baseline translation results using transformer networks for several other text-to-text sign language translation tasks.

1. Introduction

Sign Languages are the native languages of the Deaf and their main medium of communication. As visual languages, they utilize multiple complementary channels¹ to convey information [62]. This includes manual features, such as hand shape, movement and pose as well as non-manuals features, such as facial expression, mouth and movement of the head, shoulders and torso [5].

The goal of sign language translation is to either convert written language into a video of sign (production) [59, 60] or to extract an equivalent spoken language sentence from a video of someone performing continuous sign [9]. However, in the field of computer vision, much of this latter work

¹Linguists refer to these channels as articulators.

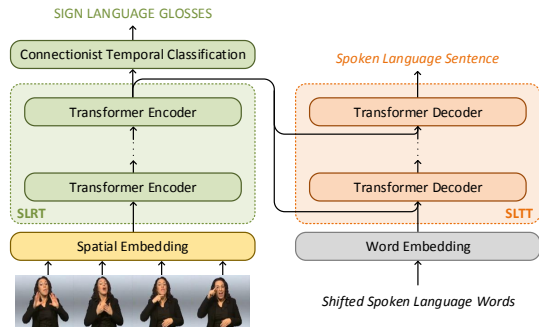


Figure 1: An overview of our end-to-end Sign Language Recognition and Translation approach using transformers.

has focused on recognising the sequence of sign glosses² (Continuous Sign Language Recognition (CSLR)) rather than the full translation to a spoken language equivalent (Sign Language Translation (SLT)). This distinction is important as the grammar of sign and spoken languages are very different. These differences include (to name a few): different word ordering, multiple channels used to convey concurrent information and the use of direction and space to convey the relationships between objects. Put simply, the mapping between speech and sign is complex and there is no simple word-to-sign mapping.

Generating spoken language sentences given sign language videos is therefore a spatio-temporal machine translation task [9]. Such a translation system requires us to accomplish several sub-tasks, which are currently unsolved:

Sign Segmentation: Firstly, the system needs to detect sign sentences, which are commonly formed using topic-comment structures [62], from continuous sign language videos. This is trivial to achieve for text based machine translation tasks [48], where the models can use punctuation marks to separate sentences. Speech-based recognition and translation systems, on the other hand, look for pauses, *e.g.* silent regions, between phonemes to segment spoken language utterances [69, 76]. There have been studies in the literature addressing automatic sign segmentation [36, 52, 55, 4, 13]. However to the best of the authors' knowledge, there is no study which utilizes sign segmentation for realizing continuous sign language translation.

²Sign glosses are spoken language words that match the meaning of signs and, linguistically, manifest as minimal lexical items.

Sign Language Recognition and Understanding: Following successful segmentation, the system needs to understand what information is being conveyed within a sign sentence. Current approaches tackle this by recognizing sign glosses and other linguistic components. Such methods can be grouped under the banner of CSLR [40, 8]. From a computer vision perspective, this is the most challenging task. Considering the input of the system is high dimensional spatio-temporal data, *i.e.* sign videos, models are required that understand what a signer looks like and how they interact and move within their 3D signing space. Moreover, the model needs to comprehend what these aspects mean in combination. This complex modelling problem is exacerbated by the asynchronous multi-articulatory nature of sign languages [51, 58]. Although there have been promising results towards CSLR, the state-of-the-art [39] can only recognize sign glosses and operate within a limited domain of discourse, namely weather forecasts [26].

Sign Language Translation: Once the information embedded in the sign sentences is understood by the system, the final step is to generate spoken language sentences. As with any other natural language, sign languages have their own unique linguistic and grammatical structures, which often do not have a one-to-one mapping to their spoken language counterparts. As such, this problem truly represents a machine translation task. Initial studies conducted by computational linguists have used text-to-text statistical machine translation models to learn the mapping between sign glosses and their spoken language translations [45]. However, glosses are simplified representations of sign languages and linguists are yet to come to a consensus on how sign languages should be annotated.

There have been few contributions towards video based continuous SLT, mainly due to the lack of suitable datasets to train such models. More recently, Camgoz *et al.* [9] released the first publicly available sign language video to spoken language translation dataset, namely PHOENIX14T. In their work, the authors proposed approaching SLT as a Neural Machine Translation (NMT) problem. Using attention-based NMT models [44, 3], they define several SLT tasks and realized the first end-to-end sign language video to spoken language sentence translation model, namely *Sign2Text*.

One of the main findings of [9] was that using gloss based mid-level representations improved the SLT performance drastically when compared to an end-to-end *Sign2Text* approach. The resulting *Sign2Gloss2Text* model first recognized glosses from continuous sign videos using a state-of-the-art CSLR method [41], which worked as a tokenization layer. The recognized sign glosses were then passed to a text-to-text attention-based NMT network [44] to generate spoken language sentences.

We hypothesize that there are two main reasons why *Sign2Gloss2Text* performs drastically better than *Sign2Text* (18.13 vs 9.58 BLEU-4 scores). Firstly, the number of sign glosses is much lower than the number of frames in the

videos they represent. By using gloss representations instead of the spatial embeddings extracted from the video frames, *Sign2Gloss2Text* avoids the long-term dependency issues, which *Sign2Text* suffers from.

We think the second and more critical reason is the lack of direct guidance for understanding sign sentences in *Sign2Text* training. Given the aforementioned complexity of the task, it might be too difficult for current Neural Sign Language Translation architectures to comprehend sign without any explicit intermediate supervision. In this paper we propose a novel Sign Language Transformer approach, which addresses this issue while avoiding the need for a two-step pipeline, where translation is solely dependent on recognition accuracy. This is achieved by jointly learning sign language recognition and translation from spatial-representations of sign language videos in an end-to-end manner. Exploiting the encoder-decoder based architecture of transformer networks [70], we propose a multi-task formalization of the joint continuous sign language recognition and translation problem.

To help our translation networks with sign language understanding and to achieve CSLR, we introduce a Sign Language Recognition Transformer (SLRT), an encoder transformer model trained using a CTC loss [2], to predict sign gloss sequences. SLRT takes spatial embeddings extracted from sign videos and learns spatio-temporal representations. These representations are then fed to the Sign Language Translation Transformer (SLTT), an autoregressive transformer decoder model, which is trained to predict one word at a time to generate the corresponding spoken language sentence. An overview of the approach can be seen in Figure 1.

The contributions of this paper can be summarized as:

- A novel multi-task formalization of CSLR and SLT which exploits the supervision power of glosses, without limiting the translation to spoken language.
- The first successful application of transformers for CSLR and SLT which achieves state-of-the-art results in both recognition and translation accuracy, vastly outperforming all comparable previous approaches.
- A broad range of new baseline results to guide future research in this field.

The rest of this paper is organized as follows: In Section 2, we survey the previous studies on SLT and the state-of-the-art in the field of NMT. In Section 3, we introduce Sign Language Transformers, a novel joint sign language recognition and translation approach which can be trained in an end-to-end manner. We share our experimental setup in Section 4. We then report quantitative results of the Sign Language Transformers in Section 5 and present new baseline results for the previously defined text-to-text translation tasks [9]. In Section 6, we share translation examples generated by our network to give the reader further qualitative insight of how our approach performs. We conclude the paper in Section 7 by discussing our findings and possible future work.

2. Related Work

Sign languages have been studied by the computer vision community for the last three decades [65, 56]. The end goal of computational sign language research is to build translation and production systems [16], that are capable of translating sign language videos to spoken language sentences and vice versa, to ease the daily lives of the Deaf [15, 6]. However, most of the research to date has mainly focused on Isolated Sign Language Recognition [35, 75, 72, 10, 63, 67], working on application specific datasets [11, 71, 23], thus limiting the applicability of such technologies. More recent work has tackled continuous data [42, 32, 17, 18], but the move from recognition to translation is still in its infancy [9].

There have been earlier attempts to realize SLT by computational linguists. However, existing work has solely focused on the text-to-text translation problem and has been very limited in size, averaging around 3000 total words [46, 57, 54]. Using statistical machine translation methods, Stein *et al.* [57] proposed a weather broadcast translation system from spoken German into German Sign Language - Deutsche Gebärdensprache (DGS) and vice versa, using the RWTH-PHOENIX-Weather-2012 (PHOENIX12) [25] dataset. Another method translated air travel information from spoken English to Irish Sign Language (ISL), spoken German to ISL, spoken English to DGS, and spoken German to DGS [45]. Ebling [22] developed an approach to translate written German train announcements into Swiss German Sign Language - Deutschschweizer Gebärdensprache (DSGS). While non-manual information has not been included in most previous systems, Ebling & Huenerfauth [24] proposed a sequence classification based model to schedule the automatic generation of non-manual features after the core machine translation step.

Conceptual video based SLT systems were introduced in the early 2000s [7]. There have been studies, such as [12], which propose recognizing signs in isolation and then constructing sentences using a language model. However, end-to-end SLT from video has not been realized until recently.

The most important obstacle to vision based SLT research has been the availability of suitable datasets. Curating and annotating continuous sign language videos with spoken language translations is a laborious task. There are datasets available from linguistic sources [53, 31] and sign language interpretations from broadcasts [14]. However, the available annotations are either weak (subtitles) or too few to build models which would work on a large domain of discourse. In addition, such datasets lack the human pose information which legacy Sign Language Recognition (SLR) methods heavily relied on.

The relationship between sign sentences and their spoken language translations are non-monotonic, as they have different ordering. Also, sign glosses and linguistic constructs do not necessarily have a one-to-one mapping with their spoken language counterparts. This made the use of available CSLR methods [42, 41] (that were designed to

learn from weakly annotated data) infeasible, as they are build on the assumption that sign language videos and corresponding annotations share the same temporal order.

To address these issues, Camgoz *et al.* [9] released the first publicly available SLT dataset, PHOENIX14T, which is an extension of the popular RWTH-PHOENIX-Weather-2014 (PHOENIX14) CSLR dataset. The authors approached the task as a spatio-temporal neural machine translation problem, which they term ‘*Neural Sign Language Translation*’. They proposed a system using Convolutional Neural Networks (CNNs) in combination with attention-based NMT methods [44, 3] to realize the first end-to-end SLT models. Following this, Ko *et al.* proposed a similar approach but used body key-point coordinates as input for their translation networks, and evaluated their method on a Korean Sign Language dataset [38].

Concurrently, there have been several advancements in the field of NMT, one of the most important being the introduction of transformer networks [70]. Transformers drastically improved the translation performance over legacy attention based encoder-decoder approaches. Also due to the fully-connected nature of the architecture, transformers are fast and easy to parallelize, which enabled them to become the new go to architecture for many machine translation tasks. In addition to NMT, transformers have achieved success in various other challenging tasks, such as language modelling [19, 77], learning sentence representations [21], multi-modal language understanding [68], activity [73] and speech recognition [34]. Inspired by their recent wide-spread success, in this work we propose a novel architecture where multiple co-dependent transformer networks are simultaneously trained to jointly solve related tasks. We then apply this architecture to the problem of simultaneous recognition and translation where joint training provides significant benefits.

3. Sign Language Transformers

In this section we introduce Sign Language Transformers which jointly learn to recognize and translate sign video sequences into sign glosses and spoken language sentences in an end-to-end manner. Our objective is to learn the conditional probabilities $p(\mathcal{G}|\mathcal{V})$ and $p(\mathcal{S}|\mathcal{V})$ of generating a sign gloss sequence $\mathcal{G} = (g_1, \dots, g_N)$ with N glosses and a spoken language sentence $\mathcal{S} = (w_1, \dots, w_U)$ with U words given a sign video $\mathcal{V} = (I_1, \dots, I_T)$ with T frames.

Modelling these conditional probabilities is a sequence-to-sequence task, and poses several challenges. In both cases, the number of tokens in the source domain is much larger than the corresponding target sequence lengths (*i.e.* $T \gg N$ and $T \gg U$). Furthermore, the mapping between sign language videos, \mathcal{V} , and spoken language sentences, \mathcal{S} , is non-monotonic, as both languages have different vocabularies, grammatical rules and orderings.

Previous sequence-to-sequence based literature on SLT can be categorized into two groups: The first group break down the problem in two stages. They consider CSLR as an

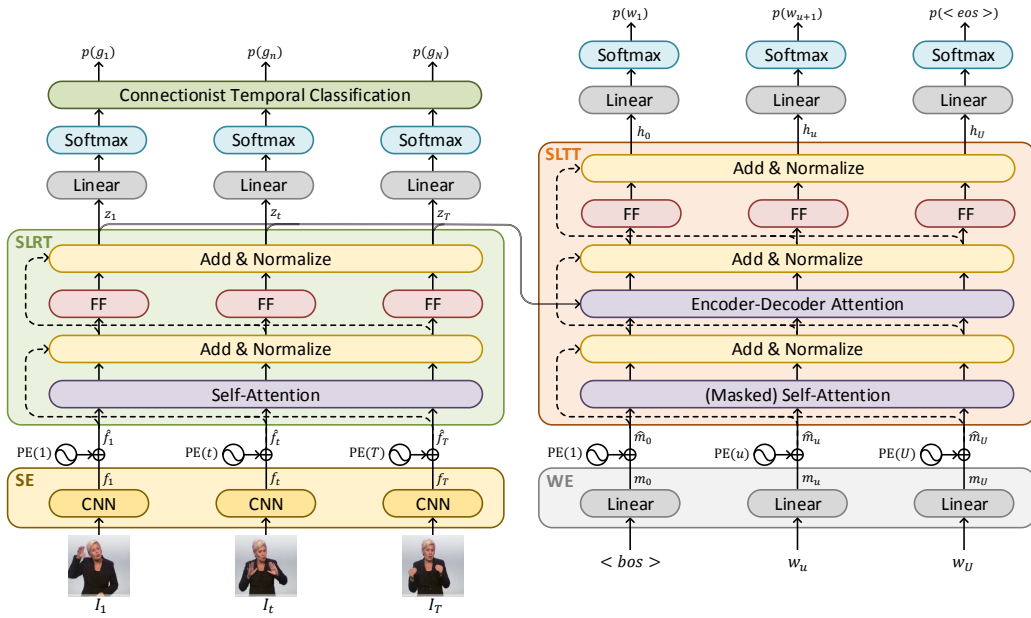


Figure 2: A detailed overview of a single layered Sign Language Transformer. (SE: Spatial Embedding, WE: Word Embedding , PE: Positional Encoding, FF: Feed Forward)

initial process and then try to solve the problem as a text-to-text translation task [12, 9]. Camgoz *et al.* utilized a state-of-the-art CSLR method [41] to obtain sign glosses, and then used an attention-based text-to-text NMT model [44] to learn the sign gloss to spoken language sentence translation, $p(\mathcal{S}|\mathcal{G})$ [9]. However, in doing so, this approach introduces an information bottleneck in the mid-level gloss representation. This limits the network’s ability to understand sign language as the translation model can only be as good as the sign gloss annotations it was trained from. There is also an inherent loss of information as a sign gloss is an incomplete annotation for linguistic study and neglects many crucial details of the vast linguistic and grammatical information present in the original sign language video.

The second group of methods focus on translation from the sign video representations to spoken language with no intermediate representation [9, 38]. These approaches attempt to learn $p(\mathcal{S}|\mathcal{V})$ directly. Given enough data and a sufficiently sophisticated network architecture, such models could theoretically realize end-to-end SLT with no need for a human-interpretable bottleneck. However, due to the lack of direct supervision guiding sign language understanding, such methods have significantly lower performance than their counterparts on the current datasets [9].

To address this, we propose to jointly learn $p(\mathcal{G}|\mathcal{V})$ and $p(\mathcal{S}|\mathcal{V})$, in an end-to-end manner. We build upon transformer networks [70] to create a unified model, which we call Sign Language Transformers (See Figure 2). We train our networks to generate spoken language sentences from sign language video representations. During training we inject intermediate gloss supervision in the form of a CTC loss into the Sign Language Recognition Transformer (SLRT) encoder. This helps our networks learn more meaningful spatio-temporal representations of the sign without

limiting the information passed to the decoder. We employ an autoregressive Sign Language Translation Transformer (SLTT) decoder which predicts one word at a time to generate the spoken language sentence translation.

3.1. Spatial and Word Embeddings

Following the classic NMT pipeline, we start by embedding our source and target tokens, namely sign language video frames and spoken language words. As word embedding we use a linear layer, which is initialized from scratch during training, to project a one-hot-vector representation of the words into a denser space. To embed video frames, we use the SpatialEmbedding approach [9], and propagate each image through CNNs. We formulate these operations as:

$$\begin{aligned} m_u &= \text{WordEmbedding}(w_u) \\ \hat{f}_t &= \text{SpatialEmbedding}(I_t) \end{aligned} \quad (1)$$

where m_u is the embedded representation of the spoken language word w_u and \hat{f}_t corresponds to the non-linear frame level spatial representation obtained from a CNN.

Unlike other sequence-to-sequence models [61, 27], transformer networks do not employ recurrence or convolutions, thus lacking the positional information within sequences. To address this issue we follow the positional encoding method proposed in [70] and add temporal ordering information to our embedded representations as:

$$\begin{aligned} \hat{f}_t &= f_t + \text{PositionalEncoding}(t) \\ \hat{m}_u &= m_u + \text{PositionalEncoding}(u) \end{aligned}$$

where PositionalEncoding is a predefined function which produces a unique vector in the form of a phase shifted sine wave for each time step.

3.2. Sign Language Recognition Transformers

The aim of SLRT is to recognize glosses from continuous sign language videos while learning meaningful spatio-temporal representations for the end goal of sign language translation. Using the positionally encoded spatial embeddings, $\hat{f}_{1:T}$, we train a transformer encoder model [70].

The inputs to SLRT are first modelled by a Self-Attention layer which learns the contextual relationship between the frame representations of a video. Outputs of the self-attention are then passed through a non-linear point-wise feed forward layer. All the operations are followed by residual connections and normalization to help training. We formulate this encoding process as:

$$z_t = \text{SLRT}(\hat{f}_t | \hat{f}_{1:T}) \quad (2)$$

where z_t denotes the spatio-temporal representation of the frame I_t , which is generated by SLRT at time step t , given the spatial representations of all of the video frames, $\hat{f}_{1:T}$.

We inject intermediate supervision to help our networks understand sign better and to guide them to learn a meaningful sign representation which helps with the main task of translation. We train the SLRT to model $p(\mathcal{G}|\mathcal{V})$ and predict sign glosses. Due to the spatio-temporal nature of the signs, glosses have a one-to-many mapping to video frames but share the same ordering.

One way to train the SLRT would be using cross-entropy loss [29] with frame level annotations. However, sign gloss annotations with such precision are rare. An alternative form of weaker supervision is to use a sequence-to-sequence learning loss functions, such as CTC [30].

Given spatio-temporal representations, $z_{1:T}$, we obtain frame level gloss probabilities, $p(g_t|\mathcal{V})$, using a linear projecting layer followed by a softmax activation. We then use CTC to compute $p(\mathcal{G}|\mathcal{V})$ by marginalizing over all possible \mathcal{V} to \mathcal{G} alignments as:

$$p(\mathcal{G}|\mathcal{V}) = \sum_{\pi \in \mathcal{B}} p(\pi|\mathcal{V}) \quad (3)$$

where π is a path and \mathcal{B} are the set of all viable paths that correspond to \mathcal{G} . We then use the $p(\mathcal{G}|\mathcal{V})$ to calculate the CSLR loss as:

$$\mathcal{L}_R = 1 - p(\mathcal{G}^*|\mathcal{V}) \quad (4)$$

where \mathcal{G}^* is the ground truth gloss sequence.

3.3. Sign Language Translation Transformers

The end goal of our approach is to generate spoken language sentences from sign video representations. We propose training an autoregressive transformer decoder model, named SLTT, which exploits the spatio-temporal representations learned by the SLRT. We start by prefixing the target spoken language sentence \mathcal{S} with the special beginning of sentence token, $\langle \text{bos} \rangle$. We then extract the positionally encoded word embeddings. These embeddings are passed to a masked self-attention layer. Although the main idea behind self-attention is the same as in the SLRT, the SLTT

utilizes a mask over the self-attention layer inputs. This ensures that each token may only use its predecessors while extracting contextual information. This masking operation is necessary, as at inference time the SLTT won't have access to the output tokens which would follow the token currently being decoded.

Representations extracted from both SLRT and SLTT self-attention layers are combined and given to an encoder-decoder attention module which learns the mapping between source and target sequences. Outputs of the encoder-decoder attention are then passed through a non-linear point-wise feed forward layer. Similar to SLRT, all the operations are followed by residual connections and normalization. We formulate this decoding process as:

$$h_{u+1} = \text{SLTT}(\hat{m}_u | \hat{m}_{1:u-1}, z_{1:T}). \quad (5)$$

SLTT learns to generate one word at a time until it produces the special end of sentence token, $\langle \text{eos} \rangle$. It is trained by decomposing the sequence level conditional probability $p(\mathcal{S}|\mathcal{V})$ into ordered conditional probabilities

$$p(\mathcal{S}|\mathcal{V}) = \prod_{u=1}^U p(w_u|h_u) \quad (6)$$

which are used to calculate the cross-entropy loss for each word as:

$$\mathcal{L}_T = 1 - \prod_{u=1}^U \sum_{d=1}^D p(\hat{w}_u^d) p(w_u^d|h_u) \quad (7)$$

where $p(\hat{w}_u^d)$ represents the ground truth probability of word w^d at decoding step u and D is the target language vocabulary size.

We train our networks by minimizing the joint loss term \mathcal{L} , which is a weighted sum of the recognition loss \mathcal{L}_R and the translation loss \mathcal{L}_T as:

$$\mathcal{L} = \lambda_R \mathcal{L}_R + \lambda_T \mathcal{L}_T \quad (8)$$

where λ_R and λ_T are hyper parameters which decides the importance of each loss function during training and are evaluated in Section 5.

4. Dataset and Translation Protocols

We evaluate our approach on the recently released PHOENIX14T dataset [9], which is a large vocabulary, continuous SLT corpus. PHOENIX14T is a translation focused extension of the PHOENIX14 corpus, which has become the primary benchmark for CSLR in recent years.

PHOENIX14T contains parallel sign language videos, gloss annotations and their translations, which makes it the only available dataset suitable for training and evaluating joint SLR and SLT techniques. The corpus includes unconstrained continuous sign language from 9 different signers with a vocabulary of 1066 different signs. Translations for these videos are provided in German spoken language with a vocabulary of 2887 different words.

The evaluation protocols on the PHOENIX14T dataset, as laid down by [9], are as follows:

Sign2Text is the end goal of SLT, where the objective is to translate directly from continuous sign videos to spoken language sentences without going via any intermediary representation, such as glosses.

Gloss2Text is a text-to-text translation problem, where the objective is to translate ground truth sign gloss sequences to German spoken language sentences. The results of these experiments act as a virtual upper bound for the available NMT translation technology. This assumption is based on the fact that perfect sign language recognition/understanding is simulated by using the ground truth gloss annotation. However, as mentioned earlier, one needs to bear in mind that gloss representations are imprecise. As glosses are textual representations of multi-channel temporal signals, they represent an information bottleneck for any translation system. This means that under ideal conditions, a *Sign2Text* system could and should outperform *Gloss2Text*. However, more sophisticated network architectures and data is needed to achieve this and hence such a goal remains a longer term objective beyond the scope of this manuscript.

Sign2Gloss2Text is the current state-of-the-art in SLT. This approach utilizes CSLR models to extract gloss sequences from sign language videos which are then used to solve the translation task as a text-to-text problem by training a *Gloss2Text* network using the CSLR predictions. *Sign2Gloss*→*Gloss2Text* is similar to *Sign2Gloss2Text* and also uses CSLR models to extract gloss sequences. However, instead of training text-to-text translation networks from scratch, *Sign2Gloss*→*Gloss2Text* models use the best performing *Gloss2Text* network, which has been trained with ground truth gloss annotations, to generate spoken language sentences from intermediate sign gloss sequences from the output of the CSLR models.

In addition to evaluating our networks in the context of the above protocols, we additionally introduce two new protocols which follow the same naming convention. *Sign2Gloss* is a protocol which essentially performs CSLR, while *Sign2(Gloss+Text)* requires joint learning of continuous sign language recognition and translation.

5. Quantitative Results

In this section we share our sign language recognition and translation experiment setups and report quantitative results. We first go over the implementation details and introduce the evaluation metrics we will be using to measure the performance of our models. We start our experiments by applying transformer networks to the text-to-text based SLT tasks, namely *Gloss2Text*, *Sign2Gloss2Text*, *Sign2Gloss*→*Gloss2Text* and report improved performance over using Recurrent Neural Network (RNN) based models. We share our *Sign2Gloss* experiments, in which we explore the effects of different types of spatial embeddings and network structures on the performance of CSLR. We then train

Sign2Text and *Sign2(Gloss+Text)* models using the best performing *Sign2Gloss* configuration and investigate the effect of different recognition loss weights on the joint recognition and translation performance. Finally, we compare our best performing models against other approaches and report state-of-the-art results.

5.1. Implementation and Evaluation Details

Framework: We used a modified version of JoeyNMT [43] to implement our Sign Language Transformers³. All components of our network were built using the PyTorch framework [50], except the CTC beam search decoding, for which we utilized the TensorFlow implementation [1].

Network Details: Our transformers are built using 512 hidden units and 8 heads in each layer. We use Xavier initialization [28] and train all of our networks from scratch. We also utilize dropout with 0.1 drop rate on transformer layers and word embeddings to mitigate over-fitting.

Performance Metrics: We use Word Error Rate (WER) for assessing our recognition models, as it is the prevalent metric for evaluating CSLR performance [40]. To measure the translation performance of our networks, we utilized BLEU [49] score (n-grams ranging from 1 to 4), which is the most common metric for machine translation.

Training: We used the Adam [37] optimizer to train our networks using a batch size of 32 with a learning rate of 10^{-3} ($\beta_1=0.9$, $\beta_2=0.998$) and a weight decay of 10^{-3} . We utilize plateau learning rate scheduling which tracks the development set performance. We evaluate our network every 100 iterations. If the development score does not decrease for 8 evaluation steps, we decrease the learning rate by a factor of 0.7. This continues until the learning rate drops below 10^{-6} .

Decoding: During the training and validation steps we employ a greedy search to decode both gloss sequences and spoken language sentences. At inference time, we utilize beam search decoding with widths ranging from 0 to 10. We also implement a length penalty [74] with α values ranging from 0 to 2. We find the best performing combination of beam width and α on the development set and use these values for the test set evaluation.

5.2. Text-to-Text Sign Language Translation

In our first set of experiments, we adapt the transformer backbone of our technique, for text-to-text sign language translation. We then evaluate the performance gain achieved over the RNN-based attention architectures.

As can be seen in Table 1, utilizing transformers for text-to-text sign language translation improved the performance across all tasks, reaching an impressive 25.35/24.54 BLEU-4 score on the development and test sets. We believe this performance gain is due to the more sophisticated attention architectures, namely self-attention modules, which learn the contextual information within both source and target sequences.

³<https://github.com/neccam/slt>

Text-to-Text Tasks (RNNs vs Transformers)	DEV					TEST				
	WER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	WER	BLEU-1	BLEU-2	BLEU-3	BLEU-4
Gloss2Text [9]	-	44.40	31.83	24.61	20.16	-	44.13	31.47	23.89	19.26
Our Gloss2Text	-	50.69	38.16	30.53	25.35	-	48.90	36.88	29.45	24.54
Sign2Gloss2Text [9]	-	42.88	30.30	23.02	18.40	-	43.29	30.39	22.82	18.13
Our Sign2Gloss2Text	-	47.73	34.82	27.11	22.11	-	48.47	35.35	27.57	22.45
Sign2Gloss→Gloss2Text [9]	-	41.08	29.10	22.16	17.86	-	41.54	29.52	22.24	17.79
Our Sign2Gloss→Gloss2Text	-	47.84	34.65	26.88	21.84	-	47.74	34.37	26.55	21.59
Video-to-Text Tasks	WER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	WER	BLEU-1	BLEU-2	BLEU-3	BLEU-4
CNN+LSTM+HMM [39]	24.50	-	-	-	-	26.50	-	-	-	-
Our Sign2Gloss	24.88	-	-	-	-	24.59	-	-	-	-
Sign2Text [9]	-	31.87	19.11	13.16	9.94	-	32.24	19.03	12.83	9.58
Our Sign2Text	-	45.54	32.60	25.30	20.69	-	45.34	32.31	24.83	20.17
Our Best Recog. Sign2(Gloss+Text)	24.61	46.56	34.03	26.83	22.12	24.49	47.20	34.46	26.75	21.80
Our Best Trans. Sign2(Gloss+Text)	24.98	47.26	34.40	27.05	22.38	26.16	46.61	33.73	26.19	21.32

Table 1: (Top) New baseline results for text-to-text tasks on Phoenix2014T [9] using transformer networks and (Bottom) Our best performing Sign Language Transformers compared against the state-of-the-art.

5.3. Sign2Gloss

To tackle the *Sign2Gloss* task, we utilize our SLRT networks. Any CNN architecture can be used as spatial embedding layers to learn the sign language video frame representation while training SLRT in an end-to-end manner. However, due to hardware limitations (graphics card memory) we utilize pretrained CNNs as our spatial embeddings. We extract frame level representations from sign videos and train our sign language transformers to learn CSLR and SLT jointly in an end-to-end manner.

In our first set of experiments, we investigate which CNN we should be using to represent our sign videos. We utilize state-of-the-art EfficientNets [66], namely B0, B4 and B7, which were trained on ImageNet [20]. We also use an Inception [64] network which was pretrained for learning sign language recognition in a CNN+LSTM+HMM setup [39]. In this set of experiments we employed a two layered transformer encoder model.

Table 2 shows that as the spatial embedding layer becomes more advanced, *i.e.* B0 vs B7, the recognition performance increases. However, our networks benefited more when we used pretrained features, as these networks had seen sign videos before and learned kernels which can embed more meaningful representations in the latent space. We then tried utilizing Batch Normalization [33] followed by a ReLU [47] to normalize our inputs and allow our networks to learn more abstract non-linear representations. This improved our results drastically, giving us a boost of nearly 7% and 6% of absolute WER reduction on the development and test sets, respectively. Considering these find-

Spatial Embedding	DEV		TEST	
	del / ins	WER	del / ins	WER
EfficientNet-B0	47.22 / 1.59	57.06	46.09 / 1.75	56.29
EfficientNet-B4	40.73 / 2.45	51.26	38.34 / 2.80	50.09
EfficientNet-B7	39.29 / 2.84	50.18	37.05 / 2.76	47.96
Pretrained CNN	21.51 / 6.10	33.90	20.29 / 5.35	33.39
+ BN & ReLU	13.54 / 5.74	26.70	13.85 / 6.43	27.62

Table 2: Impact of the Spatial Embedding Layer variants.

ings, the rest of our experiments used the batch normalized pretrained CNN features of [39] followed by ReLU.

Next, we investigated the effects of having different numbers of transformer layers. Although having a larger number of layers would allow our networks to learn more abstract representations, it also makes them prone to overfitting. To this end, we built our SLRT networks using one to six layers and evaluate their CSLR performance.

Our recognition performance initially improves with additional layers (See Table 3). However, as we continue adding more layers, our networks started to over-fit on the training data, causing performance degradation. In the light of this, for the rest of our experiments, we constructed our sign language transformers using three layers.

5.4. Sign2Text and Sign2(Gloss+Text)

In our next set of experiments we examine the performance gain achieved by unifying the recognition and translation tasks into a single model. As a baseline we trained a *Sign2Text* network, by setting our recognition loss weight λ_R to zero. We then jointly train our sign language transformers, for recognition and translation, with various weightings between the losses.

As can be seen in Table 4, jointly learning recognition and translation with equal weighting ($\lambda_R=\lambda_T=1.0$) improves the translation performance, while degrading the recognition performance compared to task specific networks. We believe this is due to scale differences of the CTC and word-level cross entropy losses. Increasing the recognition loss weight improved both the recognition and

# Layers	DEV		TEST	
	del/ins	WER	del/ins	WER
1	11.72 / 9.02	28.08	11.20 / 10.57	29.90
2	13.54 / 5.74	26.70	13.85 / 6.43	27.62
3	11.68 / 6.48	24.88	11.16 / 6.09	24.59
4	12.55 / 5.87	24.97	13.48 / 6.02	26.87
5	11.94 / 6.12	25.23	11.81 / 6.12	25.51
6	15.01 / 6.11	27.46	14.30 / 6.28	27.78

Table 3: Impact of different numbers of layers

Loss Weights		DEV		TEST	
λ_R	λ_T	WER	BLEU-4	WER	BLEU-4
1.0	0.0	24.88	-	24.59	-
0.0	1.0	-	20.69	-	20.17
1.0	1.0	35.13	21.73	33.75	21.22
2.5	1.0	26.99	22.11	27.55	21.37
5.0	1.0	24.61	22.12	24.49	21.80
10.0	1.0	24.98	22.38	26.16	21.32
20.0	1.0	25.87	20.90	25.73	20.93

Table 4: Training Sign Language Transformers to jointly learn recognition and translation with different weight on recognition loss.

the translation performance, demonstrating the value of sharing training between these related tasks.

Compared to previously published methods, our Sign Language Transformers surpass both their recognition and translation performance (See Table 1). We report a decrease of 2% WER over [39] on the test set in both *Sign2Gloss* and *Sign2(Gloss+Text)* setups. More impressively, both our *Sign2Text* and *Sign2(Gloss+Text)* networks doubled the previous state-of-the-art translation results (9.58 vs. 20.17 and 21.32 BLEU-4, respectively). Furthermore, our best performing translation *Sign2(Gloss+Text)* outperforms Camgoz *et al.*'s text-to-text based Gloss2Text translation performance (19.26 vs 21.32 BLEU-4), which was previously proposed as a pseudo upper bound on performance in [9]. This supports our claim that given more sophisticated network architectures, one would and should achieve better performance translating directly from video representations rather than doing text-to-text translation through a limited gloss representation.

6. Qualitative Results

In this section we report our qualitative results. We share the spoken language translations generated by our best performing *Sign2(Gloss+Text)* model given sign video representations (See Table 5)⁴. As the annotations in the PHOENIX14T dataset are in German, we share both the produced sentences and their translations in English.

Overall, the quality of the translations is good, and even where the exact wording differs, it conveys the same information. The most difficult translations seem to be named entities like locations which occur in limited contexts in the training data. Specific numbers are also challenging as there is no grammatical context to distinguish one from another. Despite this, the sentences produced follow standard grammar with surprisingly few exceptions.

7. Conclusion and Future Work

Sign language recognition and understanding is an essential part of the sign language translation task. Previous translation approaches heavily relied on recognition as the initial step of their system. In this paper we proposed Sign Language Transformers, a novel transformer based architecture to jointly learn sign language recognition and trans-

Reference:	im süden schwacher wind . (in the south gentle wind .)
Ours:	der wind weht im süden schwach . (the wind blows gentle in the south .)
Reference:	ähnliches wetter dann auch am donnerstag . (similar weather then also on thursday .)
Ours:	ähnliches wetter auch am donnerstag . (similar weather also on thursday .)
Reference:	ganz ähnliche temperaturen wie heute zwischen sechs und elf grad . (quite similar temperatures as today between six and eleven degrees .)
Ours:	ähnlich wie heute nacht das sechs bis elf grad . (similar as today at night that six to eleven degrees .)
Reference:	heute nacht neunzehn bis fünfzehn grad im südosten bis zwölf grad . (tonight nineteen till fifteen degrees in the southeast till twelve degrees .)
Ours:	heute nacht werte zwischen neun und fünfzehn grad im südosten bis zwölf grad . (tonight values between nine and fifteen degrees in the southeast till twelve degrees .)
Reference:	am sonntag im norden und in der mitte schauer dabei ist es im norden stürmisch . (on sunday in the north and center shower while it is stormy in the north .)
Ours:	am sonntag im norden und in der mitte niederschläge im norden ist es weiter stürmisch . (on sunday in the north and center rainfall in the north it is continuously stormy .)
Reference:	im süden und südwesten gebietsweise regen sonst recht freundlich . (in the south and southwest partly rain otherwise quite friendly .)
Ours:	im südwesten regnet es zum teil kräftig . (in the southwest partly heavy rain .)
Reference:	in der nacht sinken die temperaturen auf vierzehn bis sieben grad . (at night the temperatures lower till fourteen to seven degrees .)
Ours:	heute nacht werte zwischen sieben und sieben grad . (tonight values between seven and seven degrees .)
Reference:	heute nacht ist es meist stark bewölkt örtlich regnet oder nieselt es etwas . (tonight it is mostly cloudy locally rain or drizzle .)
Ours:	heute nacht ist es verbreitet wolkenverhangen gebietsweise regnet es kräftig . (tonight it is widespread covered with clouds partly strong rain .)
Reference:	an der saar heute nacht milde sechzehn an der elbe teilweise nur acht grad . (at the saar river tonight mild sixteen at the elbe river partly only eight degrees .)
Ours:	im rhein und südwesten macht sich morgen nur knapp über null grad . (in the rhine river and south west becomes just above zero degrees .)

Table 5: Generated spoken language translations by our Sign Language Transformers.

lation in an end-to-end manner. We utilized CTC loss to inject gloss level supervision on the transformer encoder, training it to do sign language recognition while learning meaningful representations for the end goal of sign language translation, without having an explicit gloss representation as an information bottleneck.

We evaluated our approach in the challenging PHOENIX14T dataset and report state-of-the-art sign language recognition and translation results, in some cases doubling the performance of previous translation approaches. Our first set of experiments have shown that using features which were pretrained on sign data outperformed using generic ImageNet based spatial representations. Furthermore, we have shown that jointly learning recognition and translation improved the performance across both tasks. More importantly, we have surpassed the text-to-text translation results, which was set as a virtual upper-bound, by directly translating spoken language sentences from video representations.

As future work, we would like to expand our approach to model multiple sign articulators, namely faces, hands and body, individually to encourage our networks to learn the linguistic relationship between them.

8. Acknowledgements

This work received funding from the SNSF Sinergia project ‘SMILE’ (CRSII2_160811), the European Union’s Horizon2020 research and innovation programme under grant agreement no. 762021 ‘Content4All’ and the EPSRC project ‘ExTOL’ (EP/R03298X/1). This work reflects only the author’s view and the Commission is not responsible for any use that may be made of the information it contains. We would also like to thank NVIDIA Corporation for their GPU grant.

⁴Visit *our code repository* for further qualitative examples.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: A System for Large-scale Machine Learning. In *Proceedings of the 12th Symposium on Operating Systems Design and Implementation*, 2016.
- [2] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep Audio-visual Speech Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018.
- [3] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [4] Mark Borg and Kenneth P Camilleri. Sign Language Detection “in the Wild” with Recurrent Neural Networks. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2019.
- [5] Penny Boyes-Braem and Rachel Sutton-Spence. *The Hands are the Head of the Mouth: The Mouth as Articulator in Sign Languages*. Gallaudet University Press, 2001.
- [6] Danielle Bragg, Oscar Koller, Mary Bellard, Larwan Berke, Patrick Boudreault, Annelies Braffort, Naomi Caselli, Matt Huenerfauth, Hernisa Kacorri, Tessa Verhoef, Christian Vogler, and Meredith Ringel Morris. Sign Language Recognition, Generation, and Translation: An Interdisciplinary Perspective. In *Proceedings of the International ACM SIGACCESS Conference on Computers and Accessibility (ASSETS)*, 2019.
- [7] Jan Bungeroth and Hermann Ney. Statistical Sign Language Translation. In *Proceedings of the Workshop on Representation and Processing of Sign Languages at International Conference on Language Resources and Evaluation (LREC)*, 2004.
- [8] Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, and Richard Bowden. SubUNets: End-to-end Hand Shape and Continuous Sign Language Recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [9] Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. Neural Sign Language Translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [10] Necati Cihan Camgoz, Ahmet Alp Kindiroglu, and Lale Akarun. Sign Language Recognition for Assisting the Deaf in Hospitals. In *Proceedings of the International Workshop on Human Behavior Understanding (HBU)*, 2016.
- [11] Necati Cihan Camgoz, Ahmet Alp Kindiroglu, Serpil Karabuklu, Meltem Keleşir, Ayşe Sumru Özsoy, and Lale Akarun. BosphorusSign: A Turkish Sign Language Recognition Corpus in Health and Finance Domains. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, 2016.
- [12] Xiujuan Chai, Guang Li, Yushun Lin, Zhihao Xu, Yili Tang, Xilin Chen, and Ming Zhou. Sign Language Recognition and Translation with Kinect. In *Proceedings of the International Conference on Automatic Face and Gesture Recognition (FG)*, 2013.
- [13] Neva Cherniavsky, Richard E Ladner, and Eve A Riskin. Activity Detection in Conversational Sign Language Video for Mobile Telecommunication. In *Proceedings of the International Conference on Automatic Face and Gesture Recognition (FG)*, 2008.
- [14] Helen Cooper and Richard Bowden. Learning Signs from Subtitles: A Weakly Supervised Approach to Sign Language Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [15] Helen Cooper, Brian Holt, and Richard Bowden. Sign Language Recognition. In *Visual Analysis of Humans*. Springer, 2011.
- [16] Kearsy Cormier, Neil Fox, Bencie Woll, Andrew Zisserman, Necati Cihan Camgoz, and Richard Bowden. ExTOL: Automatic recognition of British Sign Language using the BSL Corpus. In *Proceedings of the Sign Language Translation and Avatar Technology (SLTAT)*, 2019.
- [17] Runpeng Cui, Hu Liu, and Changshui Zhang. Recurrent Convolutional Neural Networks for Continuous Sign Language Recognition by Staged Optimization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [18] Runpeng Cui, Hu Liu, and Changshui Zhang. A Deep Neural Framework for Continuous Sign Language Recognition by Iterative Training. *IEEE Transactions on Multimedia*, 2019.
- [19] Zihang Dai, Zhilin Yang, Yiming Yang, William W Cohen, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive Language Models beyond a Fixed-length Context. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. ImageNet: A Large-scale Hierarchical Image Database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [21] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (ACL): Human Language Technologies*, 2019.
- [22] Sarah Ebling. *Automatic Translation from German to Synthesized Swiss German Sign Language*. PhD thesis, University of Zurich, 2016.
- [23] Sarah Ebling, Necati Cihan Camgoz, Penny Boyes Braem, Katja Tissi, Sandra Sidler-Miserez, Stephanie Stoll, Simon Hadfield, Tobias Haug, Richard Bowden, Sandrine Tornay, Marzieh Razavi, and Mathew Magimai-Doss. SMILE Swiss German Sign Language Dataset. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, 2018.
- [24] Sarah Ebling and Matt Huenerfauth. Bridging the Gap between Sign Language Machine Translation and Sign Language Animation using Sequence Classification. In *Proceedings of the 6th Workshop on Speech and Language Processing for Assistive Technologies (SPLAT)*, 2015.
- [25] Jens Forster, Christoph Schmidt, Thomas Hoyoux, Oscar Koller, Uwe Zelle, Justus H Piater, and Hermann Ney.

- RWTH-PHOENIX-Weather: A Large Vocabulary Sign Language Recognition and Translation Corpus. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, 2012.
- [26] Jens Forster, Christoph Schmidt, Oscar Koller, Martin Bellgardt, and Hermann Ney. Extensions of the Sign Language Recognition and Translation Corpus RWTH-PHOENIX-Weather. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, 2014.
- [27] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional Sequence to Sequence Learning. In *Proceedings of the ACM International Conference on Machine Learning (ICML)*, 2017.
- [28] Xavier Glorot and Yoshua Bengio. Understanding the Difficulty of Training Deep Feedforward Neural Networks. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, 2010.
- [29] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT press, 2016.
- [30] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks. In *Proceedings of the ACM International Conference on Machine Learning (ICML)*, 2006.
- [31] Thomas Hanke, Lutz König, Sven Wagner, and Silke Matthes. DGS Corpus & Dicta-Sign: The Hamburg Studio Setup. In *Proceedings of the Representation and Processing of Sign Languages: Corpora and Sign Language Technologies*, 2010.
- [32] Jie Huang, Wengang Zhou, Qilin Zhang, Houqiang Li, and Weiping Li. Video-based Sign Language Recognition without Temporal Segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [33] Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2015.
- [34] Kazuki Irie, Albert Zeyer, Ralf Schlüter, and Hermann Ney. Language Modeling with Deep Transformers. In *20th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2019.
- [35] Hamid Reza Vaezi Joze and Oscar 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.
- [36] Shujjat Khan, Donald G Bailey, and Gourab Sen Gupta. Pause detection in continuous sign language. *International Journal of Computer Applications in Technology*, 50, 2014.
- [37] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- [38] Sang-Ki Ko, Chang Jo Kim, Hyedong Jung, and Choongsang Cho. Neural Sign Language Translation based on Human Keypoint Estimation. *Applied Sciences*, 9(13), 2019.
- [39] Oscar Koller, Necati Cihan Camgoz, Richard Bowden, and Hermann Ney. Weakly Supervised Learning with Multi-Stream CNN-LSTM-HMMs to Discover Sequential Parallelism in Sign Language Videos. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.
- [40] Oscar Koller, Jens Forster, and Hermann Ney. Continuous Sign Language Recognition: Towards Large Vocabulary Statistical Recognition Systems Handling Multiple Signers. *Computer Vision and Image Understanding (CVIU)*, 141, 2015.
- [41] Oscar Koller, Sepehr Zargaran, and Hermann Ney. Re-Sign: Re-Aligned End-to-End Sequence Modelling with Deep Recurrent CNN-HMMs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [42] Oscar Koller, Sepehr Zargaran, Hermann Ney, and Richard Bowden. Deep Sign: Hybrid CNN-HMM for Continuous Sign Language Recognition. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2016.
- [43] Julia Kreutzer, Joost Bastings, and Stefan Riezler. Joey NMT: A minimalist NMT toolkit for novices. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP): System Demonstrations*, 2019.
- [44] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective Approaches to Attention-based Neural Machine Translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2015.
- [45] Sara Morrissey. *Data-driven Machine Translation for Sign Languages*. PhD thesis, Dublin City University, 2008.
- [46] Sara Morrissey, Harold Somers, Robert Smith, Shane Gilchrist, and Sandipan Dandapat. Building a Sign Language Corpus for Use in Machine Translation. In *Proceedings of the Representation and Processing of Sign Languages: Corpora and Sign Language Technologies*, 2010.
- [47] Vinod Nair and Geoffrey E Hinton. Rectified Linear Units Improve Restricted Boltzmann Machines. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010.
- [48] Graham Neubig. Neural Machine Translation and Sequence-to-Sequence Models: A Tutorial. *arXiv:1703.01619*, 2017.
- [49] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: A Method for Automatic Evaluation of Machine Translation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2002.
- [50] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic Differentiation in PyTorch. In *Proceedings of the Advances in Neural Information Processing Systems Workshops (NIPSW)*, 2017.
- [51] Wendy Sandler. Sign Language and Modularity. *Lingua*, 89(4), 1993.
- [52] Pinar Santemiz, Oya Aran, Murat Saraclar, and Lale Akarun. Automatic Sign Segmentation from Continuous Signing via Multiple Sequence Alignment. In *Proceedings of the IEEE International Conference on Computer Vision Workshops (ICCVW)*, 2009.
- [53] Adam Schembri, Jordan Fenlon, Ramas Rentelis, Sally Reynolds, and Kearsy Cormier. Building the British Sign Language Corpus. *Language Documentation & Conservation (LD&C)*, 7, 2013.
- [54] Christoph Schmidt, Oscar Koller, Hermann Ney, Thomas Hoyoux, and Justus Piater. Using Viseme Recognition to Improve a Sign Language Translation System. In *Proceedings of the International Workshop on Spoken Language Translation*, 2013.

- [55] Frank M Shipman, Satyakiran Duggina, Caio DD Monteiro, and Ricardo Gutierrez-Osuna. Speed-Accuracy Tradeoffs for Detecting Sign Language Content in Video Sharing Sites. In *Proceedings of the International ACM SIGACCESS Conference on Computers and Accessibility (ASSETS)*, 2017.
- [56] Thad Starner, Joshua Weaver, and Alex Pentland. Real-time American Sign Language Recognition using Desk and Wearable Computer Based Video. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 20(12), 1998.
- [57] Daniel Stein, Christoph Schmidt, and Hermann Ney. Analysis, Preparation, and Optimization of Statistical Sign Language Machine Translation. *Machine Translation*, 26(4), 2012.
- [58] William C Stokoe. Sign Language Structure. *Annual Review of Anthropology*, 9(1), 1980.
- [59] Stephanie Stoll, Necati Cihan Camgoz, Simon Hadfield, and Richard Bowden. Sign Language Production using Neural Machine Translation and Generative Adversarial Networks. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2018.
- [60] Stephanie Stoll, Necati Cihan Camgoz, Simon Hadfield, and Richard Bowden. Text2Sign: Towards Sign Language Production using Neural Machine Translation and Generative Adversarial Networks. *International Journal of Computer Vision (IJCV)*, 2020.
- [61] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to Sequence Learning with Neural Networks. In *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [62] Rachel Sutton-Spence and Bencie Woll. *The Linguistics of British Sign Language: An Introduction*. Cambridge University Press, 1999.
- [63] Muhammed Mirac Suzgun, Hilal Ozdemir, Necati Cihan Camgoz, Ahmet Kindiroglu, Dogac Basaran, Cengiz Togay, and Lale Akarun. Hospisign: An Interactive Sign Language Platform for Hearing Impaired. In *Proceedings of the International Conference on Computer Graphics, Animation and Gaming Technologies (Eurasia Graphics)*, 2015.
- [64] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence.*, 2017.
- [65] Shinichi Tamura and Shingo Kawasaki. Recognition of Sign Language Motion Images. *Pattern Recognition*, 21(4), 1988.
- [66] Mingxing Tan and Quoc V Le. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- [67] Sandrine Tornay, Marzieh Razavi, Necati Cihan Camgoz, Richard Bowden, and Mathew Magimai-Doss. HMM-based Approaches to Model Multichannel Information in Sign Language Inspired from Articulatory Features-based Speech Processing. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2019.
- [68] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal Transformer for Unaligned Multimodal Language Sequences. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [69] Jan P van Hemert. Automatic Segmentation of Speech. *IEEE Transactions on Signal Processing*, 39(4), 1991.
- [70] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All You Need. In *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, 2017.
- [71] Hanjie Wang, Xiujuan Chai, and Xilin Chen. Sparse Observation (SO) Alignment for Sign Language Recognition. *Neurocomputing*, 175, 2016.
- [72] Hanjie Wang, Xiujuan Chai, Xiaopeng Hong, Guoying Zhao, and Xilin Chen. Isolated Sign Language Recognition with Grassmann Covariance Matrices. *ACM Transactions on Accessible Computing*, 8(4), 2016.
- [73] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7794–7803, 2018.
- [74] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s Neural Machine Translation System: Bridging the Gap Between Human and Machine Translation. *arXiv:1609.08144*, 2016.
- [75] Fang Yin, Xiujuan Chai, and Xilin Chen. Iterative Reference Driven Metric Learning for Signer Independent Isolated Sign Language Recognition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [76] Norimasa Yoshida. *Automatic Utterance Segmentation in Spontaneous Speech*. PhD thesis, Massachusetts Institute of Technology, 2002.
- [77] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.