# Signing at Scale: Learning to Co-Articulate Signs for Large-Scale Photo-Realistic Sign Language Production

Ben Saunders, Necati Cihan Camgoz, Richard Bowden
University of Surrey
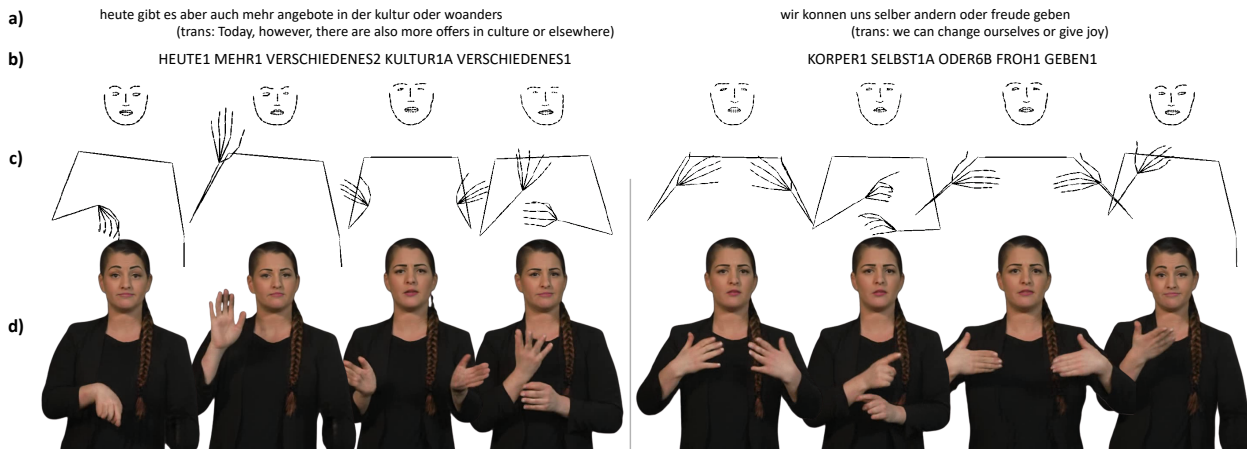{b.saunders, n.camgoz, r.bowden}@surrey.ac.uk

Figure 1. **Photo-Realistic Sign Language Production:** Given a spoken language sentence from an unconstrained domain of discourse (a), an initial translation is conducted to a gloss sequence (b). FS-NET next produces a co-articulated continuous skeleton pose sequence from dictionary signs (c), which SIGNGAN generates into a photo-realistic sign language video in a given style (d).

## Abstract

*Sign languages are visual languages, with vocabularies as rich as their spoken language counterparts. However, current deep-learning based Sign Language Production (SLP) models produce under-articulated skeleton pose sequences from constrained vocabularies and this limits applicability. To be understandable and accepted by the deaf, an automatic SLP system must be able to generate co-articulated photo-realistic signing sequences for large domains of discourse.*

*In this work, we tackle large-scale SLP by learning to co-articulate between dictionary signs, a method capable of producing smooth signing while scaling to unconstrained domains of discourse. To learn sign co-articulation, we propose a novel Frame Selection Network (FS-NET) that improves the temporal alignment of interpolated dictionary signs to continuous signing sequences. Additionally, we propose SIGNGAN, a pose-conditioned human synthesis model that produces photo-realistic sign language videos direct from skeleton pose. We propose a novel keypoint-based loss function which improves the quality of synthesized hand images.*

*We evaluate our SLP model on the large-scale meineDGS (mDGS) corpus, conducting extensive user evaluation showing our FS-NET approach improves co-articulation of interpolated dictionary signs. Additionally, we show that SIGNGAN significantly outperforms all baseline methods for quantitative metrics, human perceptual studies and native deaf signer comprehension.*

## 1. Introduction

Sign languages are rich visual languages with large lexical vocabularies [53] and intricate co-articulated movements of both manual (hands and body) and non-manual (facial) features. Sign Language Production (SLP), the automatic translation from spoken language sentences to sign language sequences, must be able to produce photo-realistic continuous signing for large domains of discourse to be useful to the deaf communities.

Prior deep-learning approaches to SLP have either pro-

duced concatenated isolated sequences that disregard the natural co-articulation between signs [54, 69] or continuous sequences end-to-end [23, 45, 47, 69] which suffer from under-articulation [44]. Furthermore, these methods have struggled to generalise beyond the limited domain of weather [15].

In this paper, we propose an SLP method to produce photo-realistic continuous sign language videos direct from unconstrained spoken language sequences. Firstly, we translate from spoken language to gloss[1] sequences. We next learn the temporal co-articulation between gloss-based dictionary signs, modelling the temporal prosody of sign language [3].

To model sign co-articulation, we propose a novel Frame Selection Network (FS-NET) that learns the optimal subset of frames that best represents a continuous signing sequence (Fig. 2 middle). We build a transformer encoder with cross-attention [59] to predict a temporal alignment path supervised by Dynamic Time Warping (DTW).

The resulting skeleton pose sequences are subsequently used to condition a video-to-video synthesis model capable of generating photo-realistic sign language videos, named SIGNGAN (Fig. 2 right). Due to the natural presence of motion blur in sign language datasets from fast moving hands [16], a classical application of a hand discriminator leads to an increase in blurred hand generation. To avoid this, we propose a novel keypoint-based loss that significantly improves the quality of hand image synthesis in our photo-realistic signer generation module. To enable training on diverse sign language datasets, we propose a method for controllable video generation that models a multi-modal distribution of sign language videos in different styles.

Our deep-learning based SLP model is able to generalise to large domains of discourse, as it is trivial to increase vocabulary with a few examples of this new sign in a continuous signing context. We conduct extensive deaf user evaluation on a translation protocol of mDGS [21], showing that FS-NET improves the natural signing motion of interpolated dictionary sequences and is overwhelmingly preferred to baseline SLP methods [48]. Additionally, we achieve state-of-the-art back translation performance on RWTH-PHOENIX-Weather-2014**T** (PHOENIX14**T**) with a 43% improvement over baselines, highlighting the understandable nature of our approach.

Furthermore, we evaluate SIGNGAN using the high quality Content4All (C4A) dataset [5], outperforming state-of-the-art synthesis methods [7, 55, 63, 64] for quantitative evaluation and human perception studies. Finally, we conduct a further deaf user evaluation to show that SIGNGAN is more understandable than the skeletal sequences previously used to represent sign [45].

The contributions of this paper can be summarised as:

- The first SLP model to produce large-scale sign language sequences from an unconstrained domain of discourse to a level understandable by a native deaf signer

- A novel Frame Selection Network, FS-NET, that learns to co-articulate between dictionary signs via a monotonic alignment to continuous sequences

- A method to generate photo-realistic continuous sign language videos, SIGNGAN, with a novel hand keypoint loss that improves the hand synthesis quality

- Extensive user evaluation of our proposed approach, showing preference of our proposed method, alongside state-of-the-art back translation results

## 2. Related Work

**Sign Language Production** The initial focus of computational sign language technology was Sign Language Recognition (SLR) [12, 19, 29] with few works tackling unconstrained SLR [9, 26, 30]. More recently, focus has shifted to Sign Language Translation (SLT) [4, 8, 28].

Sign Language Production (SLP), the translation from spoken to sign language, has been historically tackled using animated avatars [10, 27, 37] with rules-based co-articulation that does not generalise to unseen sequences [50].

Initial deep learning-based SLP methods concatenated isolated signs with no regards for natural co-articulation [54, 69]. Recently, continuous SLP methods have directly regressed sequences of multiple signs [23, 45, 47–49], but exhibit under-articulated signing motion due to regression to the mean. To overcome under-articulation, we avoid generating pose directly and learn the optimal temporal alignment between dictionary and continuous sign sequences.

In addition, prior work has represented sign language as skeleton pose sequences, which have been shown to reduce the deaf comprehension compared to a photo-realistic production [60]. Previous works have attempted photo-realistic signer generation [11, 46, 55], but of low-resolution isolated signs. In this work, we produce high-resolution photo-realistic continuous sign language videos directly from spoken language input, from unrestricted domains of discourse.

**Pose-Conditioned Human Synthesis** Generative Adversarial Networks (GANs) [18] have achieved impressive results in image [24, 42, 64, 71] and, more recently, video generation tasks [36, 58, 61–63]. Specific to pose-conditioned human synthesis, there has been concurrent research focusing on the generation of whole body [1, 35, 38, 51, 56, 72], face [13, 32, 68] and hand [34, 57, 67] images.
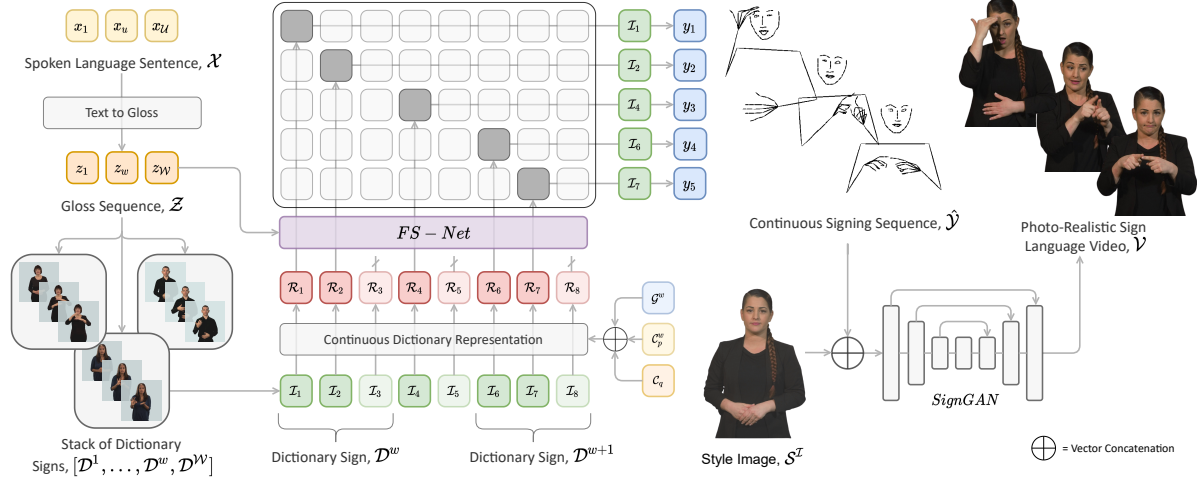
---

[1]Glosses are a written representation of sign that follow sign language ordering and grammar, defined as minimal lexical items [53].

Figure 2. Overview of our proposed large-scale SLP method. An initial Text to Gloss (left) animates an interpolated dictionary sequence, $\mathcal{I}$, with a Frame Selection Network (FS-NET), learning the temporal alignment (middle) to a continuous signing sequence, $\mathcal{Y}$. Finally, SIGNGAN generates a photo-realistic sign language video, $\mathcal{V}$, (right) from the continuous skeleton pose and a given style image, $\mathcal{S}^{\mathcal{I}}$.

However, there has been no research into accurate hand generation in the context of full body synthesis, with current methods failing to generate high-quality hand images [60]. Due to the hands being high fidelity objects, they are often overlooked in model optimisation. Chan *et al.* introduced FaceGAN for high resolution face generation [7], but no similar work has been proposed for the more challenging task of hand synthesis in the context of sign language, where hand to hand interaction is ubiquitous. In this work, we propose a keypoint-based loss to enhance hand synthesis.

The task of human motion transfer, transferring motion from source to target videos via keypoint extraction, is relevant to our task [7, 66, 70]. However, there has been limited research into the generation of novel poses, which we produce from a given spoken language sentence. Additionally, works have attempted to produce unseen appearances in a few-shot manner [62, 68], but continue to produce only a single style at inference.

**Sign Language Co-Articulation** Sign language co-articulation can be defined as "the articulatory influence of one phonetic element on another across more than one intervening element" [20] and is an important distinction between isolated and natural continuous signing [40].

Co-articulation involves both the motion and duration of signs, with a particular focus on the transition between signs [40]. The boundaries of a sign are also modified depending on the context, with continuous signing typically produced faster than their isolated counterparts [50]. In this work, we model temporal co-articulation by learning the optimal alignment between isolated signs and continuous signing sequences, predicting the duration, boundary and transition of each sign in context.

## 3. Large-Scale Photo-Realistic SLP

The true aim of a large-scale SLP model is to translate a spoken language sequence from an unconstrained domain of discourse, $\mathcal{X} = (x_1, ..., x_{\mathcal{U}})$ with $\mathcal{U}$ words, to a continuous photo-realistic sign language video, $\mathcal{V}^{\mathcal{S}} = (v_1, ..., v_{\mathcal{T}})$ with $\mathcal{T}$ frame. This is a challenging task due to the large vocabulary of unconstrained signing and the intricate spatial nature of sign, with a requirement for temporal co-articulation indicative of natural continuous signing.

We approach this problem as a multi-stage sequence-to-sequence task. Firstly, spoken language is translated to sign gloss, $\mathcal{Z} = (z_1, ..., z_{\mathcal{W}})$, as an intermediate representation (Sec. 3.1). Next, our FS-NET model co-articulates between gloss-based dictionary signs to produce a full continuous signing sequence, $\mathcal{Y} = (y_1, ..., y_{\mathcal{T}})$ (Sec. 3.2). Finally, given $\mathcal{Y}$ and a style image, $\mathcal{S}^{\mathcal{I}}$, our video-to-video signer generation module generates a photo-realistic sign language video, $\mathcal{Z}^{\mathcal{S}}$ (Sec. 3.3). An overview of our approach can be seen in Fig. 2. In the remainder of this section, we shall describe each component of our approach in detail.

### 3.1. Text to Gloss

Given a spoken language sequence, $\mathcal{X}$, we first translate to a sign language grammar and order, represented by a gloss sequence, $\mathcal{Z} = (z_1, ..., z_{\mathcal{W}})$ with $\mathcal{W}$ glosses (Fig. 2 left). We formulate this as a sequence-to-sequence problem, due to the non-monotonic relationship between the two sequences of different lengths. We use an encoder-decoder transformer [59] to perform this translation, formalised as:

$$f_t = E_{T2G}(x_t | x_{1:\mathcal{T}}) \tag{1}$$

$$g_{w+1} = D_{T2G}(g_w|g_{1:w-1}, f_{1:\mathcal{T}}) \tag{2}$$

where $f_t$ and $g_w$ are the encoded source and target tokens respectively and $g_0$ is the encoding of the special $< \text{bos} >$ token. The output gloss tokens can be computed as $z_w = \text{argmax}_i(g_w)$ until the special $< \text{eos} >$ token is predicted.

## 3.2. Gloss to Pose

Next, motivated by the monotonic relationship between glosses and signs, we produce a continuous signing pose sequence, $\hat{\mathcal{Y}} = (y_1, ..., y_{\mathcal{T}})$ with $\mathcal{T}$ frames, from the translated gloss sequence, $\mathcal{Z}$, using a learnt co-articulation of dictionary signs. We first encode the gloss sequence using a transformer encoder with self-attention, as:

$$h_w = E_{G2S}(z_w|z_{1:\mathcal{W}}) \tag{3}$$

where $h_w$ is the encoded gloss token for step $w$. We next collect a dictionary sample, $\mathcal{D}^w$, for every sign present in the gloss vocabulary. By definition, dictionary signs contain accurate and articulated sign content. Furthermore, it is trivial to expand to larger domains of discourse, simply collecting dictionary examples of the expanded vocabulary.

**Interpolated Dictionary Representation** Given the translated gloss sequence, $\mathcal{Z}$, we create a stack of ordered dictionary signs, $[\mathcal{D}^1, ..., \mathcal{D}^w, \mathcal{D}^{\mathcal{W}}]$ (Bottom left of Fig. 2). As in previous works [47], we represent each dictionary sign as a sequence of skeleton pose, $\mathcal{D}^w = (s_1^w, ..., s_{\mathcal{P}^w}^w)$ with $\mathcal{P}^w$ frames. We first convert the stack of dictionary signs into a continuous sequence by linearly interpolating between neighbouring signs for a predefined fixed $\mathcal{N}_{LI}$ frames. The final interpolated dictionary sequence, $\mathcal{I} = (\mathcal{I}_1, ..., \mathcal{I}_{\mathcal{Q}})$ with $\mathcal{Q}$ frames, is the combination of skeleton pose and the respective linear interpolation.

We next build a continuous dictionary sequence representation to be used as input to FS-NET. Alongside the skeleton pose of $\mathcal{I}$, we learn a gloss embedding, $\mathcal{G}^w$, unique to each gloss in the vocabulary, with a separate shared embedding for all interpolation frames, $\mathcal{G}^{LI}$. Additionally, we use a counter embedding proposed by Saunders *et al.* [45], expanded to both a specific counter, $\mathcal{C}_p^w$, relating to the progression of each dictionary sign and a global counter, $\mathcal{C}_q$, relating to the progress of the full sequence, $\mathcal{I}$. The final continuous dictionary representation, $\mathcal{R} = (\mathcal{R}_1, ..., \mathcal{R}_{\mathcal{Q}})$ with $\mathcal{Q}$ frames, is constructed by concatenating the corresponding skeleton, gloss and counter embeddings per frame, as:

$$\mathcal{R}_q = [s_p^w, \mathcal{G}^w, \mathcal{C}_p^w, \mathcal{C}_q] \tag{4}$$

where frame $q$ represents a time step $p$ frames into gloss $w$.

**Frame Selection Network** To co-articulate between dictionary signs, we propose a Frame Selection Network (FS-NET) that learns to predict the temporal alignment to a continuous signing sequence, $\mathcal{Y} = (y_0, ..., y_{\mathcal{T}})$ with $\mathcal{T}$ frames (Fig. 2 middle). We note that this is a monotonic sequence-to-sequence task, due to the matching order of signing and the different sequence lengths ($\mathcal{Q} \neq \mathcal{T}$).

Formally, FS-NET predicts a discrete sparse monotonic temporal alignment path, $\hat{\mathcal{A}} \in \mathbb{R}^{\mathcal{Q} \times \mathcal{Q}}$:

$$\hat{\mathcal{A}} = \text{FS-NET}(\mathcal{R}, h_{1:\mathcal{W}}) \tag{5}$$

where $\hat{\mathcal{A}}$ contains binary decisions representing either frame selection or skipping. Fig. 2 shows an example alignment that skips the production of frames $3, 5$ and $8$ in the output sequence, removing redundant frames to create a smoother co-articulated continuous sequence. We build FS-NET as a transformer encoder [59] with an additional cross-attention to the encoded gloss sequence. To produce the co-articulated continuous signing pose sequence, $\hat{\mathcal{Y}}$, a matrix multiplication can be applied between $\mathcal{I}$ and $\hat{\mathcal{A}}$, as:

$$\hat{y} = \mathcal{I} \times \hat{\mathcal{A}} \tag{6}$$

This enables the mapping between varied length sequences, with the end of sequence prediction determined as the alignment selection of the final dictionary frame.

**Dynamic Time Warping Supervision** In practice, directly predicting the 2D alignment, $\hat{\mathcal{A}}$, provides weak gradients due to the sparse nature of the alignment. We therefore propose to train FS-NET using a Dynamic Time Warping (DTW) supervision signal [2] designed to learn the optimal monotonic temporal alignment. We pre-compute the DTW path, $\mathcal{A}^* = \text{DTW}(\mathcal{Q}, \mathcal{T})$, between the interpolated dictionary sequence, $\mathcal{I}$, and the target continuous sequence, $\mathcal{Y}$. Due to the intractability of 2D alignment path prediction, we collapse the alignment down to a 1D sequence during training, $\hat{\text{A}} \in \mathbb{R}^{\mathcal{Q}} = \text{argmax}_q(\hat{\mathcal{A}})$. This enables a temporal mask prediction over $\mathcal{I}$, selecting which frames of the interpolated dictionary sequence to animate in turn to create a continuous sequence.

We argue that for the majority of sequences (88% for mDGS), $\mathcal{Q} >> \mathcal{T}$, due to the faster tempo of continuous sign [40]. We thus assume that no frames are added during temporal alignment, only removed. To train FS-NET, we compute a cross entropy loss $\mathcal{L}_{CE}$ between the predicted 1D temporal alignment, $\hat{\text{A}} \in \mathbb{R}^{\mathcal{Q}}$, and the ground truth DTW alignment, $\text{A}^* \in \mathbb{R}^{\mathcal{Q} \times 1}$, as:

$$\mathcal{L}_{CE}(\hat{\text{A}}, \text{A}^*) = -\frac{1}{\mathcal{Q}} \sum_{q=1}^{\mathcal{Q}} \text{A}_q^* \cdot \log(\hat{\text{A}}_q) \tag{7}$$

The final continuous sign pose sequence, $\hat{\mathcal{Y}} = (y_1, ..., y_{\mathcal{T}})$, is produced as shown in Eq. 6.

## 3.3. Pose to Video

To generate a photo-realistic sign language video, $\mathcal{V}^{\mathcal{S}}$, conditioned on the produced sign pose sequence, $\hat{\mathcal{Y}}$, we propose a method for video-to-video signer generation, SIGN-GAN (Fig. 2 right). Taking inspiration from [7], in the conditional GAN setup, a generator network, $G$, competes in a min-max game against a multi-scale discriminator, $D = (D_1, D_2, D_3)$. The goal of $G$ is to synthesise images of similar quality to ground-truth images, in order to fool $D$. Conversely, the aim of $D$ is to discern the "fake" images from the "real" images. For our purposes, $G$ synthesises images of a signer, $v^{\mathcal{S}}$ given a human pose, $y_t$, and a style image, $\mathcal{S}^{\mathcal{I}}$.

Following [24], we introduce skip connections to the architecture of $G$ in a *U-Net* structure [43] between each down-sampling layer $i$ and up-sampling layer $n-i$, where $n$ is the total number of up-sampling layers. Skip connections propagate pose information across the networks, enabling the generation of fine-grained details. Specifically, we add skip connections between each down-sampling layer $i$ and up-sampling layer $n - i$, where $n$ is the total number of up-sampling layers.

**Controllable Video Generation** To enable training on diverse sign language datasets, we use a style-controllable video generation approach [46]. A style image, $\mathcal{S}^{\mathcal{I}}$, is provided to condition synthesis alongside the pose sequence, as seen in Figure 2. SIGNGAN learns to associate the given style, $\mathcal{S}$, with the person-specific aspects of the corresponding target image, $v_t^{\mathcal{S}}$, such as the clothing or face, but disentangle the signer-invariant skeleton pose.

Controllable generation allows SIGNGAN to make use of the variability in signer appearance in the data. A multi-modal distribution of sign language videos in different styles, $\mathcal{V}^{S}$, can be produced, where $S \in \{1, N_S\}$ represents the styles seen during training [2].

**Hand Keypoint Loss** Previous pose-conditioned human synthesis methods have failed to generate realistic and accurate hand images [60]. To enhance the quality of hand synthesis, we introduce a novel loss that operates in the keypoint space, as shown in Figure 3. A pre-trained 2D hand pose estimator [17], $H$, is used to extract hand keypoints, $k_H$, from cropped hand regions (*i.e.* a 60x60 patch centered around the middle knuckle), $v_H$, as $k_H = H(v_H)$. We avoid operating in the image space due to the existence of blurry hand images in the dataset, whereas the extracted keypoints are invariant to motion-blur. A hand keypoint discriminator, $D_H$, then attempts to discern between the
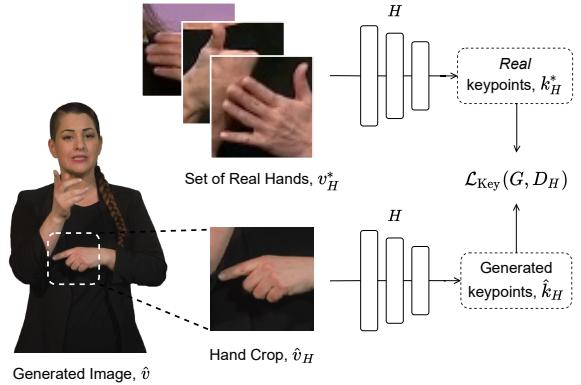


Figure 3. Hand keypoint loss overview. A keypoint discriminator, $D_H$, compares keypoints extracted from generated hands, $\hat{k}_H$, and real hands, $k_H^\star$.

"real" keypoints, $k_H^\star = H(v_H)$, and the "fake" keypoints, $\hat{k}_H = H(G(y_H))$, leading to the objective:

$$\mathcal{L}_{\text{KEY}}(G, D_H) = \mathbb{E}_{y_H, z_H}[\log D_H(k_H^\star)]$$
$$+ \mathbb{E}_{y_H}[\log(1 - D_H(\hat{k}_H))] \tag{8}$$

**Full Objective** In standard image-to-image translation frameworks [24,64], $G$ is trained using a combination of adversarial and perceptual losses. We update the multi-scale adversarial loss, $\mathcal{L}_{GAN}(G, D)$, to reflect the controllable generation with a joint conditioning on sign pose, $y_t$, and style image, $\mathcal{S}^{\mathcal{I}}$:

$$\mathcal{L}_{GAN}(G, D) = \sum_{i=1}^{k} \mathbb{E}_{y_t, z_t}[\log D_i(z_t \mid y_t, \mathcal{S}^{\mathcal{I}})]$$
$$+ \mathbb{E}_{y_t}[\log(1 - D_i(G(y_t, \mathcal{S}^{\mathcal{I}}) \mid y_t, \mathcal{S}^{\mathcal{I}}))] \tag{9}$$

where $k = 3$ reflects the multi-scale discriminator. The adversarial loss is supplemented with two feature-matching losses; $\mathcal{L}_{FM}(G, D)$, the discriminator feature-matching loss presented in pix2pixHD [64], and $\mathcal{L}_{VGG}(G, D)$, the perceptual reconstruction loss [25] which compares pre-trained VGGNet [52] features at multiple layers of the network. Our full SIGNGAN objective, $\mathcal{L}_{Total}$, is a weighted sum of these, alongside our proposed hand keypoint loss (Eq. 8), as:

$$\mathcal{L}_{Total} = \min_{G}((\max_{D_i} \sum_{i=1}^{k} \mathcal{L}_{GAN}(G, D_i))$$

$$+ \lambda_{FM} \sum_{i=1}^{k} \mathcal{L}_{FM}(G, D_i) + \lambda_{VGG} \mathcal{L}_{VGG}(G(y_t, I^S), z_t)$$

$$+ \lambda_{\text{KEY}} \mathcal{L}_{\text{KEY}}(G, D_H)) \tag{10}$$

where $k = 3$ and $\lambda_{FM}, \lambda_{VGG}, \lambda_{\text{KEY}}$ weight the contributions of each loss.

---

[2]For qualitative examples (e.g. in Fig 4), we share a single signer appearance, as we have consent from this signer to use their appearance for publication purposes.
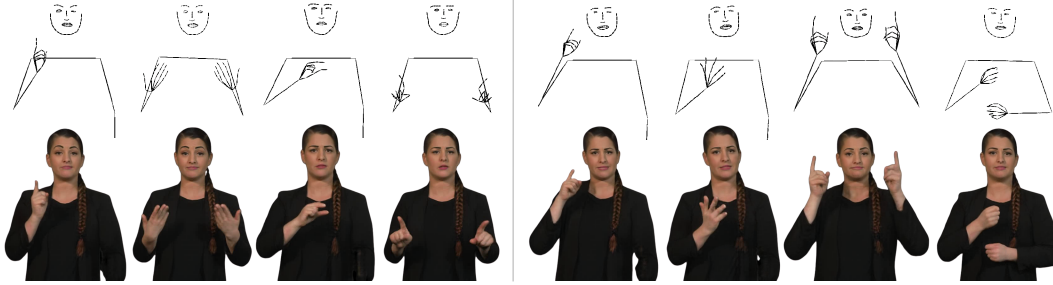
Figure 4. Example photo-realistic frames with skeleton pose produced using FS-NET and photo-realistic video generated using SIGNGAN.

## 4. Experiments

In this section, we evaluate our large-scale photo-realistic SLP approach. We outline our experimental setup then perform quantitative, qualitative & user evaluation.

### 4.1. Experimental Setup

To train our large-scale SLP approach, we set a new translation protocol on the Meine DGS (mDGS) corpus[3] [21], a large German Sign Language - Deutsche Gebärdensprache (DGS) linguistic resource capturing free-form signing from 330 deaf participants, with a vocabulary of 10,042 glosses. To adapt the corpus for translation, we segment the free-flowing discourse into 40,230 segments of German sentences, sign gloss translations and sign language videos. We pre-process the mDGS gloss annotations [31] and create two protocols, with either gloss variants included (mDGS-V) or removed (mDGS). We publicly release these translation protocols[4] to facilitate the future growth in large-scale SLP and SLT research, with further details provided in the appendix. A license must be obtained from the University of Hamburg to use mDGS for computational research.

For additional experiments, we use the benchmark PHOENIX14**T** dataset [4] from the constrained weather broadcast domain, with setup and skeletal pose configuration as in [45]. We collect exhaustive dictionary examples of every DGS sign present in mDGS and PHOENIX14**T**, trimmed to remove the sign onset and offset. For samples without expressive mouthings, we insert the facial features present in a example of the respective gloss from the continuous training set. For photo-realistic video generation, we use the C4A dataset [5] due to its high video quality and diverse interpreter appearance. We use a heat-map representation as pose condition, with each skeletal limb plotted on a separate feature channel.

We build our Text to Gloss models with 2 layers, 4 heads and hidden size of 128 and our FS-NET with 2 layers, 4 heads and 64 hidden size. We set the interpolation frames, $\mathcal{N}_{LI}$, to 5 and the learning rate to $10^{-3}$. Our code is based on JoeyNMT [33], and implemented using PyTorch [41].

---

[3]With permission from the University of Hamburg.

[4]https://github.com/BenSaunders27/meineDGS-Translation-Protocols

### 4.2. Quantitative Evaluation

#### 4.2.1 Text to Gloss

We first evaluate our Text to Gloss translation described in Sec. 3.1. Table 1 shows a performance of 21.93 BLEU-4 on PHOENIX14**T**, outperforming [45] (20.23) but falling short of [39] (23.17) who use larger training data. Translation performance is considerably lower on both meineDGS-Variants (mDGS-V) and mDGS due to the larger domain, showing that further research is required to scale the task to larger vocabularies.

| Dataset: | DEV SET | | TEST SET | |
|---|---|---|---|---|
| | BLEU-4 | ROUGE | BLEU-4 | ROUGE |
| mDGS-V | 1.96 | 24.51 | 1.16 | 25.34 |
| mDGS | 3.17 | 32.93 | 3.08 | 32.52 |
| PHOENIX14**T** | **21.93** | **57.25** | **20.08** | **56.63** |

Table 1. Text to Gloss results on mDGS and PHOENIX14**T**.

#### 4.2.2 Gloss to Pose

**Back Translation** Back translation has developed as the state-of-the-art SLP evaluation metric [45]. We train an SLT model [8] on PHOENIX14**T** with skeleton pose sequences generated using our FS-NET approach. Table 2 shows considerable performance gains (43%) compared to baseline methods on the Gloss to Pose task [44,45,47,48]. This highlights the increased comprehension provided by FS-NET compared to baseline end-to-end regression methods, with an ability to overcome the poor quality of the PHOENIX14**T** dataset. Furthermore, it can be seen that interpolated dictionary sequences (Comparable to that of

| Approach: | DEV SET | | TEST SET | |
|---|---|---|---|---|
| | BLEU-4 | ROUGE | BLEU-4 | ROUGE |
| Progressive Transformers [45] | 11.93 | 34.01 | 10.43 | 32.02 |
| Adversarial Training [44] | 13.16 | 36.75 | 12.16 | 34.19 |
| Mixture Density Networks [47] | 13.14 | 39.06 | 11.94 | 35.19 |
| Mixture of Motion Primitives [48] | 13.32 | 37.58 | 12.67 | 35.61 |
| Interpolated Dictionary Sequence | 16.28 | 38.11 | 16.27 | 36.95 |
| **FS-NET (Ours)** | **19.14** | **40.94** | **18.78** | **40.60** |

Table 2. Back translation results on the PHOENIX14**T** dataset for the *Gloss to Pose* task.

| Approach: | DEV SET | | TEST SET | |
|---|---|---|---|---|
| | BLEU-4 | ROUGE | BLEU-4 | ROUGE |
| Progressive Transformers [45] | 11.82 | 33.18 | 10.51 | 32.46 |
| Adversarial Training [44] | 12.65 | 33.68 | 10.81 | 32.74 |
| Mixture Density Networks [47] | 11.54 | 33.40 | 11.68 | 33.19 |
| Mixture of Motion Primitives [48] | 14.03 | **37.76** | 13.30 | 36.77 |
| **FS-NET (Ours)** | **16.92** | 35.74 | **21.10** | **42.57** |

Table 3. Back translation results on the PHOENIX14**T** dataset for the *Text to Pose* task.

[54]) achieve worse back translation results, highlighting the effect of FS-NET co-articulation for comprehension.

Additionally, Table 3 shows further state-of-the-art results on the full pipeline of text to pose, with an initial text to gloss translation and a subsequent sign animation. This highlights the effectiveness of the full spoken language to photo-realistic video pipeline required for true SLP.

**Sign User Evaluation** We next perform extensive user evaluation with native signers, animating our skeleton pose outputs using SIGNGAN. All baselines are also generated by SIGNGAN, to alleviate visual differences in comparison. In total, 10 participants completed our sign user evaluation, of which all were fluent signers and 20% were deaf. We provide all the generated user evaluation videos in the supplementary materials.

We first compare the comprehension of FS-NET compared to the state-of-the-art deep learning based SLP method [48]. We show participants pairs of generated videos from the same sequence, asking to select which video was more understandable. Table 4 shows how our productions were unanimously preferred to baselines for both mDGS and PHOENIX14**T**. This overwhelming result highlights both the increased comprehension of FS-NET, alongside the inability of previous methods to scale to unconstrained domains of discourse.

We next evaluate how understandable our large-scale sign productions are in isolation. We show each participant a produced video alongside a list of 10 signs, of which 5 are signed in the video, and ask them to select which signs they believe are being signed. For FS-NET productions, an average of 4.8 signs were recognised for each video. This shows our productions are easily understandable by native signers, an essential result for accurate large-scale SLP.

Our final user evaluation evaluates how co-articulated

| Dataset | FS-NET | Baseline [48] | Equal |
|---|---|---|---|
| mDGS | **95%** | 0% | 5% |
| PHOENIX14**T** | **95%** | 0% | 5% |

Table 4. Comprehension user evaluation results, showing the percentage of participants who chose productions from FS-NET or a baseline [48] to be more understandable, or equal.

| | FS-NET | Isolated | Equal |
|---|---|---|---|
| Non-Trimmed | **100**% | 0% | 0% |
| Trimmed | 40% | **47**% | 13% |

Table 5. Co-articulation user evaluation results, showing the percentage of participants who believed the video with the smoothest transitions was from FS-NET, isolated concatenation or equal.

our FS-NET productions are. We show participants two videos of the same sequence, one isolated dictionary sequence and one co-articulated continuous video generated by FS-NET, and ask them to select which they believed had the most smooth transitions between signs. We first evaluate dictionary signs without trimming the sign onset and offset, with Table 5 showing our productions were unanimously preferred. Moving to trimmed dictionary signs, FS-NET productions were preferred 40% of the time, with 13% equal preference. This highlights the effectiveness of FS-NET at improving co-articulation between dictionary signs and temporally aligning to continuous signing sequences.

### 4.2.3 Pose to Video

Finally, we evaluate our photo-realistic sign language video approach, SIGNGAN. We compare the performance of SIGNGAN against state-of-the-art image-to-image and video-to-video translation methods [7, 55, 63, 64], conditioned on skeletal pose images. We measure the quality of synthesized images using the following metrics; 1) SSIM: Structural Similarity [65] over the full image. 2) Hand SSIM: SSIM metric over a crop of each hand. 3) Hand Pose: Absolute distance between 2D hand keypoints of the produced and ground truth hand images, using a pre-trained hand pose estimation model [17]. 4) FID: Fréchet Inception Distance [22] over the full image.

**Baseline Comparison** We first compare SIGNGAN to baseline methods for photo-realistic generation given a sequence of ground truth poses as input. Table 6 shows results on the C4A data, with SIGNGAN outperforming all baselines particularly for the Hand SSIM and FID. We believe this is due to the improved quality of synthesized hand images by using the proposed hand keypoint loss.

| | SSIM ↑ | Hand SSIM ↑ | Hand Pose ↓ | FID ↓ |
|---|---|---|---|---|
| EDN [7] | 0.737 | 0.553 | 23.09 | 41.54 |
| vid2vid [63] | 0.750 | 0.570 | 22.51 | 56.17 |
| Pix2PixHD [64] | 0.737 | 0.553 | 23.06 | 42.57 |
| Stoll *et al.* [55] | 0.727 | 0.533 | 23.17 | 64.01 |
| SIGNGAN (Ours) | **0.759** | **0.605** | **22.05** | **27.75** |

Table 6. Baseline model comparison results of photo-realistic sign language video generation.

|  | SSIM $\uparrow$ | Hand SSIM $\uparrow$ | Hand Pose $\downarrow$ | FID $\downarrow$ |
|---|---|---|---|---|
| Baseline | 0.743 | 0.582 | 22.87 | 39.33 |
| Hand Discriminator | 0.738 | 0.565 | 22.81 | 39.22 |
| Hand Keypoint Loss | **0.759** | **0.605** | **22.05** | **27.75** |

Table 7. Ablation study results of SignGAN

**Ablation Study** We perform an ablation study of Sign-GAN, with results in Table 7. As suggested in Sec. 3.3, the hand discriminator performs poorly for both SSIM and hand SSIM, due to the generation of blurred hands. However, our proposed hand keypoint loss increases model performance considerably and particularly for Hand SSIM, emphasizing the importance of an adversarial loss invariant to blurring.

**Perceptual Study** We perform an additional perceptual study of our photo-realistic generation, showing participants pairs of 10 second videos generated by SignGAN and a corresponding baseline method. Participants were asked to select which video was more visually realistic, with a separate focus on the body and hands. In total, 46 participant completed the study, of which 28% were signers, each viewing 2 randomly selected videos from each baseline. Table 8 shows the percentage of participants who preferred the outputs of SignGAN to the baseline method. It can be seen that SignGAN outputs were unanimously preferred for both body (96.2% average) and hand (95.6% average) synthesis. Vid2vid [63] was the strongest contender, with our productions preferred only 85% of the time.

**Deaf User Evaluation** Our final user evaluation compares the comprehension of photo-realistic videos against the previously-used skeletal pose representation [45]. We provided 5 30-second videos of ground-truth skeletal sequences and corresponding photo-realistic videos to deaf participants, asking them to rate each video out of 5 for understandability. Synthesised videos were rated higher for comprehension, at 3.9 compared to 3.2 for skeletal sequences. This suggests that photo-realistic content is more understandable to a deaf signer than a skeleton sequence.

### 4.3. Qualitative Evaluation

We show example generated photo-realistic frames in Fig. 4, highlighting the production quality. We provide further qualitative evaluation in supplementary materials.

|  | Body | Hand |
|---|---|---|
| EDN [7] | 100% | 97.8% |
| vid2vid [63] | 85.9% | 84.8% |
| Pix2PixHD [64] | 98.9% | 100% |
| Stoll *et al.* [55] | 100% | 100% |

Table 8. Perceptual study results, showing the percentage of participants who preferred SignGAN to the baseline model.

## 5. Potential Negative Societal Impact

We acknowledge the potential use of SLP technology to remove the reliance on human interpreters. However, we see this work as enabling a larger provision of signed content, especially where interpretation doesn't exist [14]. We also recognise the potential harm if this technology produced incorrect sign language content, particularly in emergency settings. Although this paper significantly advances the SLP field, we would like to state that SLP technology is still under development and should not yet be relied upon.

## 6. Conclusion

Large-scale photo-realistic SLP is important to provide high quality signing content to deaf communities. In this paper, we proposed the first SLP method to achieve both large-scale signing and photo-realistic video generation. We proposed FS-Net, which learns to co-articulate between dictionary signs by modelling the optimal temporal alignment to continuous sequences. Furthermore, we proposed SignGAN to produce photo-realistic sign language videos. We proposed a novel keypoint-based loss function that improves the quality of hand synthesis, operating in the keypoint space to avoid issues caused by motion blur.

We showed how our approach can scale to unconstrained domains of discourse and be understood by native signers, with considerable state-of-the-art PHOENIX14**T** back translation performance. Additionally, we performed extensive user evaluation showing our approach increases the realism of interpolated dictionary signs, can be understood by native deaf signers and is overwhelmingly preferred to baseline methods. Finally, we showed that SignGAN outperforms all baseline methods for quantitative metrics, human evaluation and native deaf signer comprehension.

Our approach is limited by the current performance of text to gloss translation for large-scale domains. Available gloss annotations are limited, making sign language translation tasks a low-resource machine translation task [39]. Improvements on both architectures and datasets are required to compete with spoken language Neural Machine Translation (NMT) methods. For future work, we plan to tackle spatial co-articulation between dictionary signs.

# References

[1] Guha Balakrishnan, Amy Zhao, Adrian V Dalca, Fredo Durand, and John Guttag. Synthesizing Images of Humans in Unseen Poses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[2] Donald J Berndt and James Clifford. Using Dynamic Time Warping to Find Patterns in Time Series. In *AAA1-94 Workshop on Knowledge Discovery in Databases*, 1994. 4

[3] Diane Brentari, Joshua Falk, Anastasia Giannakidou, Annika Herrmann, Elisabeth Volk, and Markus Steinbach. Production and Comprehension of Prosodic Markers in Sign Language Imperatives. *Frontiers in Psychology*, 2018. 2

[4] Necati Cihan Camgöz, 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. 2, 6

[5] Necati Cihan Camgöz, Ben Saunders, Guillaume Rochette, Marco Giovanelli, Giacomo Inches, Robin Nachtrab-Ribback, and Richard Bowden. Content4All Open Research Sign Language Translation Datasets. In *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, 2021. 2, 6

[6] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 12

[7] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody Dance Now. In *Proceedings of the IEEE International Conference on Computer Vision (CVPR)*, 2019. 2, 3, 5, 7, 8

[8] Necati Cihan Camgöz, Oscar Koller, Simon Hadfield, and Richard Bowden. Sign Language Transformers: Joint End-to-end Sign Language Recognition and Translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 6

[9] Helen Cooper and Richard Bowden. Large Lexicon Detection of Sign Language. In *International Workshop on Human-Computer Interaction*, 2007. 2

[10] Stephen Cox, Michael Lincoln, Judy Tryggvason, Melanie Nakisa, Mark Wells, Marcus Tutt, and Sanja Abbott. TESSA, a System to Aid Communication with Deaf People. In *Proceedings of the ACM International Conference on Assistive Technologies*, 2002. 2

[11] Runpeng Cui, Zhong Cao, Weishen Pan, Changshui Zhang, and Jianqiang Wang. Deep Gesture Video Generation With Learning on Regions of Interest. *IEEE Transactions on Multimedia*, 2019. 2

[12] 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. 2

[13] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and Controllable Face Image Generation via 3D Imitative-Contrastive Learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[14] Jules Dickinson. *Sign Language Interpreting in the Workplace*. Gallaudet University Press, 2017. 8

[15] 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. 2

[16] 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. 2

[17] Liuhao Ge, Zhou Ren, Yuncheng Li, Zehao Xue, Yingying Wang, Jianfei Cai, and Junsong Yuan. 3D Hand Shape and Pose Estimation from a Single RGB Image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 5, 7

[18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, 2014. 2

[19] Kirsti Grobel and Marcell Assan. Isolated Sign Language Recognition using Hidden Markov Models. In *IEEE International Conference on Systems, Man, and Cybernetics*, 1997. 2

[20] Michael Andrew Grosvald. *Long-Distance Coarticulation: A Production and Perception Study of English and American Sign Language*. University of California, Davis, 2009. 3

[21] Thomas Hanke, Lutz König, Sven Wagner, and Silke Matthes. DGS Corpus & Dicta-Sign: The Hamburg Studio Setup. In *4th Workshop on the Representation and Processing of Sign Languages: Corpora and Sign Language Technologies (CSLT 2010), Valletta, Malta*, 2010. 2, 6, 12

[22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, 2017. 7

[23] Wencan Huang, Wenwen Pan, Zhou Zhao, and Qi Tian. Towards Fast and High-Quality Sign Language Production. In *Proceedings of the 29th ACM International Conference on Multimedia*, 2021. 2

[24] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-Image Translation with Conditional Adversarial Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 5

[25] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 5

[26] Timor Kadir, Richard Bowden, Eng-Jon Ong, and Andrew Zisserman. Minimal Training, Large Lexicon, Uncon-

strained Sign Language Recognition. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2004. 2

[27] Kostas Karpouzis, George Caridakis, S-E Fotinea, and Eleni Efthimiou. Educational Resources and Implementation of a Greek Sign Language Synthesis Architecture. *Computers & Education (CAEO)*, 2007. 2

[28] Sang-Ki Ko, Chang Jo Kim, Hyedong Jung, and Choongsang Cho. Neural Sign Language Translation based on Human Keypoint Estimation. *Applied Sciences*, 2019. 2

[29] Oscar Koller. Quantitative Survey of the State of the Art in Sign Language Recognition. *arXiv preprint arXiv:2008.09918*, 2020. 2

[30] 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)*, 2015. 2

[31] Reiner Konrad, Thomas Hanke, Gabriele Langer, Susanne König, Lutz König, Rie Nishio, and Anja Regen. Public DGS Corpus: Annotation Conventions. Technical report, Project Note AP03-2018-01, DGS-Korpus project, IDGS, Hamburg University, 2018. 6, 12

[32] Marek Kowalski, Stephan J Garbin, Virginia Estellers, Tadas Baltrušaitis, Matthew Johnson, and Jamie Shotton. CON-FIG: Controllable Neural Face Image Generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 2

[33] 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)*, 2019. 6

[34] Yahui Liu, Marco De Nadai, Gloria Zen, Nicu Sebe, and Bruno Lepri. Gesture-to-Gesture Translation in the Wild via Category-Independent Conditional Maps. In *Proceedings of the 27th ACM International Conference on Multimedia*, 2019. 2

[35] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose Guided Person Image Generation. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2

[36] Arun Mallya, Ting-Chun Wang, Karan Sapra, and Ming-Yu Liu. World-Consistent Video-to-Video Synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 2

[37] John McDonald, Rosalee Wolfe, Jerry Schnepp, Julie Hochgesang, Diana Gorman Jamrozik, Marie Stumbo, Larwan Berke, Melissa Bialek, and Farah Thomas. Automated Technique for Real-Time Production of Lifelike Animations of American Sign Language. *Universal Access in the Information Society (UAIS)*, 2016. 2

[38] Yifang Men, Yiming Mao, Yuning Jiang, Wei-Ying Ma, and Zhouhui Lian. Controllable Person Image Synthesis with Attribute-Decomposed GAN. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[39] Amit Moryossef, Kayo Yin, Graham Neubig, and Yoav Goldberg. Data Augmentation for Sign Language Gloss Translation. In *Proceedings of the Biennial Machine Translation Summit, International Workshop on Automatic Translation for Signed and Spoken Languages*, 2021. 6, 8

[40] Lucie Naert, Caroline Larboulette, and Sylvie Gibet. Coarticulation Analysis for Sign Language Synthesis. In *International Conference on Universal Access in Human-Computer Interaction*, 2017. 3, 4

[41] 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 *NIPS Autodiff Workshop*, 2017. 6

[42] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. *arXiv preprint arXiv:1511.06434*, 2015. 2

[43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional Networks for Biomedical Image Segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MIC-CAI))*, 2015. 5

[44] Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Adversarial Training for Multi-Channel Sign Language Production. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2020. 2, 6, 7

[45] Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Progressive Transformers for End-to-End Sign Language Production. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 2, 4, 6, 7, 8

[46] Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. AnonySign: Novel Human Appearance Synthesis for Sign Language Video Anonymisation. *arXiv preprint arXiv:2107.10685*, 2021. 2, 5

[47] Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Continuous 3D Multi-Channel Sign Language Production via Progressive Transformers and Mixture Density Networks. 2021. 2, 4, 6, 7

[48] Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Mixed SIGNals: Sign Language Production via a Mixture of Motion Primitives. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 2, 6, 7

[49] Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. Skeletal Graph Self-Attention: Embedding a Skeleton Inductive Bias into Sign Language Production. *arXiv preprint arXiv:2112.05277*, 2021. 2

[50] Jérémie Segouat. A Study of Sign Language Coarticulation. *ACM SIGACCESS Accessibility and Computing*, 2009. 2, 3

[51] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuiliere, and Nicu Sebe. Deformable GANs for Pose-Based Human Image Generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[52] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5

[53] William C Stokoe. Sign Language Structure. *Annual Review of Anthropology*, 1980. 1, 2

[54] Stephanie Stoll, Necati Cihan Camgöz, 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. 2, 7

[55] Stephanie Stoll, Necati Cihan Camgöz, 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. 2, 7, 8

[56] Hao Tang, Song Bai, Li Zhang, Philip HS Torr, and Nicu Sebe. XingGAN for Person Image Generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 2

[57] Hao Tang, Wei Wang, Dan Xu, Yan Yan, and Nicu Sebe. GestureGAN for Hand Gesture-to-Gesture Translation in the wild. In *Proceedings of the 26th ACM International Conference on Multimedia*, 2018. 2

[58] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing Motion and Content for Video Generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2, 3, 4

[60] Lucas Ventura, Amanda Duarte, and Xavier Giró-i Nieto. Can Everybody Sign Now? Exploring Sign Language Video Generation from 2D Poses. In *ECCV Sign Language Recognition, Translation and Production Workshop*, 2020. 2, 3, 5

[61] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating Videos with Scene Dynamics. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. 2

[62] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. Few-shot Video-to-Video Synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2, 3

[63] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-Video Synthesis. In *Advances in Neural Information Processing Systems (NIPS)*, 2018. 2, 7, 8

[64] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 5, 7, 8

[65] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image Quality Assessment: From Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing*, 2004. 7

[66] Dongxu Wei, Xiaowei Xu, Haibin Shen, and Kejie Huang. GAC-GAN: A General Method for Appearance-Controllable Human Video Motion Transfer. *IEEE Transactions on Multimedia*, 2020. 3

[67] Zhenyu Wu, Duc Hoang, Shih-Yao Lin, Yusheng Xie, Liangjian Chen, Yen-Yu Lin, Zhangyang Wang, and Wei Fan. MM-Hand: 3D-Aware Multi-Modal Guided Hand Generative Network for 3D Hand Pose Synthesis. *arXiv preprint arXiv:2010.01158*, 2020. 2

[68] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-Shot Adversarial Learning of Realistic Neural Talking Head Models. In *Proceedings of the IEEE International Conference on Computer Vision (CVPR)*, 2019. 2, 3

[69] Jan Zelinka and Jakub Kanis. Neural Sign Language Synthesis: Words Are Our Glosses. In *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020. 2

[70] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara Berg. Dance Dance Generation: Motion Transfer for Internet Videos. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019. 3

[71] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017. 2

[72] Zhen Zhu, Tengteng Huang, Baoguang Shi, Miao Yu, Bofei Wang, and Xiang Bai. Progressive Pose Attention Transfer for Person Image Generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

# Appendices

## meineDGS (mDGS) Translation Protocol

In this appendix, we provide further details of our released translation protocols on the meineDGS (mDGS) dataset [21]. The public mDGS linguistic corpus can be accessed at https://www.sign-lang.uni-hamburg.de/meinedgs/, containing 330 sequences of free-flowing discourse between two deaf participants, with each around 10 minutes in length. Additionally, detailed spoken language transcripts, frame-level gloss annotations and 2D pose estimation sequences [6] are provided. Discourse is centered around a wide variety of topics, age groups and format, with further details available on the mDGS website.

To adapt the mDGS corpus for use as a translation dataset, we segment the free-flowing discourse data into 40,230 segments of German sentences, sign gloss translations and respective sign language videos. Sequence segmentation was performed using spoken language sentence boundaries, with corresponding frame boundaries provided. The title of each segment (e.g. 1583882A-X) contains the title of the original discourse sequence as given in the *Transcript* column (e.g. 1583882), the corresponding participant camera (A or B) and the position of the extracted segment in the original discourse sequence (a numerical value X).

Table 9 and 10 show detailed statistics of the mDGS-V and mDGS protocols, respectively. Gloss variants used in mDGS-V give distinction between sign variants, with each containing the same meaning but with differing motion. We chose to retain these variants to provide more challenging baselines for the community. Further public annotation conventions are outlined in [31], which we follow. Additionally, gloss frame alignments are provided as *GLOSS/start-frame/stop-frame* (e.g. BUCHSTABE1/11/34). The translation protocols are publicly available at https://github.com/BenSaunders27/meineDGS-Translation-Protocols, detailing *filename*, *camera*, *ger_text*, *gloss*, *start_time* and *stop_time*.

|  | Sign Gloss | | | German | | |
|---|---|---|---|---|---|---|
|  | Train | Dev | Test | Train | Dev | Test |
| segments | 40,230 | 4,996 | 4,977 | ←————— same | | |
| frames | 6,146,153 | 764,451 | 758,883 | ←————— same | | |
| vocab. | 10,042 | 4,644 | 4,620 | 18,680 | 6,224 | 6,231 |
| tot. words | 215,392 | 26,855 | 26,969 | 389,427 | 48,376 | 48,551 |
| tot. OOVs | - | 371 | 339 | - | 1,103 | 1,171 |
| singletons | 2,681 | - | - | 8,909 | - | - |

Table 9. Key statistics of the meineDGS-Variants (mDGS-V) dataset split.

|  | Sign Gloss | | | German | | |
|---|---|---|---|---|---|---|
|  | Train | Dev | Test | Train | Dev | Test |
| segments | 40,230 | 4,996 | 4,977 | ←————— same | | |
| frames | 6,146,153 | 764,451 | 758,883 | ←————— same | | |
| vocab. | 4,337 | 2,490 | 2,487 | 18,680 | 6,224 | 6,231 |
| tot. words | 215,392 | 26,855 | 26,969 | 389,427 | 48,376 | 48,551 |
| tot. OOVs | - | 118 | 112 | - | 1,103 | 1,171 |
| singletons | 778 | - | - | 8,909 | - | - |

Table 10. Key statistics of the meineDGS (mDGS) dataset split.

To use the mDGS dataset for computational research, a licence must be obtained from the University of Hamburg[5]. Release of these protocols does not imply permission for use or provide a license. Written permission is required from the dataset owner. Please adhere to the data ownership policies and ensure you have the correct rights of use.

---

[5] https://www.sign-lang.uni-hamburg.de/meinedgs/