

How2Sign: A Large-scale Multimodal Dataset for Continuous American Sign Language

Amanda Duarte^{1,2*}
Kenneth DeHaan⁶

Shruti Palaskar⁴
Florian Metze^{4,5}

Lucas Ventura¹
Jordi Torres^{1,2}

Deepti Ghadiyaram⁵
Xavier Giro-i-Nieto^{1,2,3*}

¹Universitat Politècnica de Catalunya ²Barcelona Supercomputing Center ³Institut de Robòtica i Informàtica Industrial, CSIC-UPC
⁴Carnegie Mellon University ⁵Facebook AI ⁶Gallaudet University

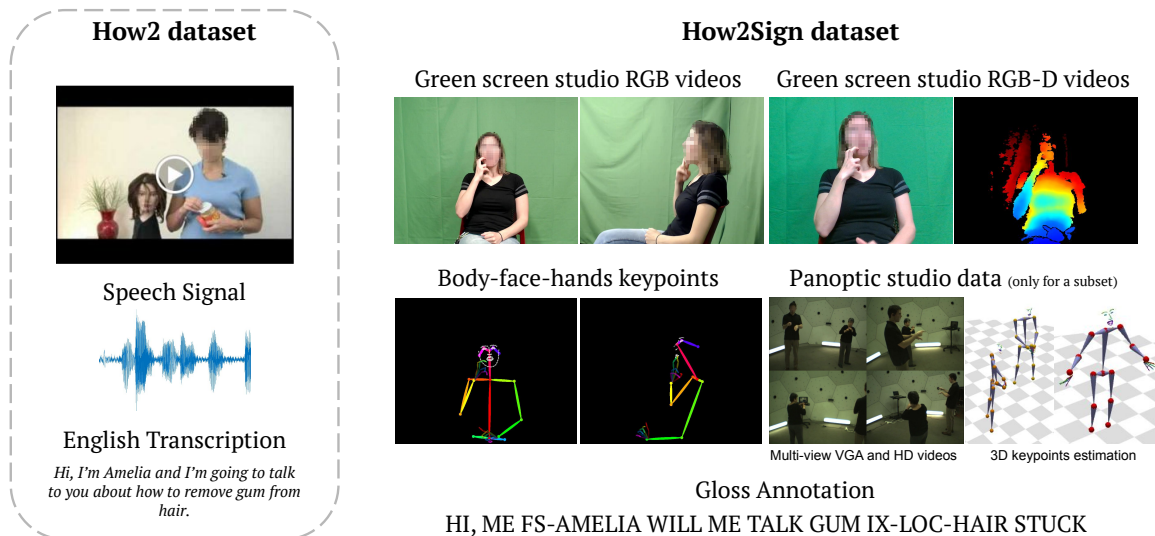


Figure 1: The **How2Sign** dataset consists of over 80 hours of multiview sign language videos and aligned modalities.

Abstract

One of the factors that have hindered progress in the areas of sign language recognition, translation, and production is the absence of large annotated datasets. Towards this end, we introduce How2Sign, a multimodal and multiview continuous American Sign Language (ASL) dataset, consisting of a parallel corpus of more than 80 hours of sign language videos and a set of corresponding modalities including speech, English transcripts, and depth. A three-hour subset was further recorded in the Panoptic studio enabling detailed 3D pose estimation. To evaluate the potential of How2Sign for real-world impact, we conduct a study with ASL signers and show that synthesized videos using our dataset can indeed be understood. The study further gives insights on challenges that computer vision should address in order to make progress in this field.

Dataset website: <http://how2sign.github.io/>

*Corresponding authors: {amanda.duarte,xavier.giro}@upc.edu

1. Introduction

Sign languages (SL) are the primary means of communication for an estimated 466 million deaf¹ or hard-of-hearing people worldwide [1]. Like any other natural language, sign languages are consistently evolving and have structure directed by a set of linguistic rules [3]. They differ from spoken languages and do not have standard written forms, e.g. American Sign Language (ASL) is not a sign form of English. Although sign languages are used by millions of people everyday to communicate, the vast majority of communications technologies nowadays are designed to support spoken or written language, but not sign languages. At the same time, most hearing people do not know a sign language; as a result, many communication barriers exist for deaf sign language users [6, 7, 14].

¹We follow the recognized convention of using the upper-cased word Deaf which refers to the culture and describes members of the community of sign language users and the lower-cased word deaf describes the hearing status[37].

Promising recent works in sign language processing² [12, 30, 33, 41, 40, 19] have shown that modern computer vision and machine learning architectures can help break down these barriers for sign language users. Improving such models could make technologies that are primarily designed for non-sign language users, *e.g.* voice-activated services, text-based systems, spoken-media based content, *etc.*, more accessible to the Deaf community. Other possibilities include automatic transcription of signed content, which would help facilitating the communication between sign and non-sign language users, as well as real-time interpreting when human interpreters are not available, and many other educational tools and applications [6].

However, training such models requires large amounts of data. The availability of public large-scale datasets suitable for machine learning is very limited, especially when it comes to *continuous sign language* datasets, *i.e.*, where the data needs to be segmented and annotated at the sentence level. Currently, there is no ASL dataset large enough to be used with recent deep learning approaches.

In order to instigate the advance in the area of research that involves sign language processing, in this paper we introduce the *How2Sign* dataset. *How2Sign* is a large-scale collection of multimodal and multiview sign language videos in American Sign Language (ASL) for over 2500 instructional videos selected from the existing *How2* dataset [27]. Figure 1 shows samples of the data contained in the dataset. Working in close collaboration with native ASL signers and professional interpreters, we collected more than *80 hours* of multi-view and multimodal (recorded with multiple RGB and a depth sensor) ASL videos, and corresponding gloss annotations [22]. In addition, a three-hour subset was further recorded at the Panoptic studio [17], a geodesic dome setup equipped with hundreds of cameras and sensors, which enables detailed 3D reconstruction and pose estimation. This subset paves the way for vision systems to understand the 3D geometry of sign language.

Our contributions can be summarized as follows: a) We present *How2Sign*, a large-scale multimodal and multiview continuous American Sign Language dataset that consists of more than *80 hours of American Sign Language videos*, with sentence-level alignment for more than 35k sentences. It features a vocabulary of 16k English words that represent more than two thousand instructional videos from a broad range of categories; b) Our dataset comes with a rich set of annotations including gloss, category labels, as well automatically extracted 2D keypoints for more than 6M frames. What is more, a subset of the dataset was re-recorded in the Panoptic studio with more than 500 cameras that enabled high quality 3D keypoints estimation for around 3 hours of

videos; c) We conduct a study with ASL signers that showed that videos generated using our dataset can be understood to a certain extent, and at the same time gave insights on challenges that the research community can address in this field.

2. Background and Related Work

In this section we discuss some of the challenges that comes with sign languages that can be interesting to the computer vision community, as well as an overview of the current publicly available sign language datasets.

2.1. Sign Language

Sign languages are visual languages that use two types of features to convey information: *manual* that includes handshape, palm orientation, movement and location and; *non-manual markers* that are movement of the head (nod/shake/tilt), mouth (mouthing), eyebrows, cheeks, facial grammar (or facial expressions) and eye gaze [32]. All these features need to be taken into account while recognizing, translating or generating signs in order to capture the complete meaning of the sign. This makes sign language processing a challenging set of tasks for computer vision.

When it comes to continuous sign language, a simple concatenation of isolated signs is not enough to correctly recognize, translate or generate a complete sentence and neglects the underlying rich grammatical and linguistic structures of sign language that differ from spoken language. Besides the fact that the alignment between sign and spoken language sequences are usually unknown and non monotonic [12], the transitions between signs must also be taken into account. Usually, the beginning of a sign is modified depending on the previous sign, and the end of the same sign is modified depending on the following sign making them visually different in the isolated and continuous scenarios [3]. This phenomenon is called “co-articulation” and brings an extra challenge for tasks with continuous sign language [2].

2.2. Sign Language datasets

One of the most important factors that has hindered the progress of sign language processing research is the absence of large scale annotated datasets [6]. Many existing sign language datasets contain isolated signs [10, 4, 18, 21, 23, 34]. Such data may be important for certain scenarios (*e.g.*, creating a dictionary, or as a resource for those who are learning a sign language), but most real-world use cases of sign language processing involve natural conversational with complete sentences (*i.e.* *continuous sign language*).

A number of continuous sign language datasets have been collected over the years mainly for linguistic purposes. *SIGNUM* [35] and the *BSL Corpus* [31] were recorded in

²For brevity, we follow [6] and use the term *sign language processing* to refer to the set of sign language recognition, translation and production tasks.

Name	Language	Vocab.	Duration (h)	Signers	Modalities					
					Multiview	Transcription	Gloss	Pose	Depth	Speech
Video-Based CSL [16]	CSL	178	100	50	✗	✓	✗	✓	✓	✗
SIGNUM [35]	DGS	450	55	25	✗	✓	✓	✗	✗	✗
RWTH-Phoenix-2014T [12]	DGS	3k	11	9	✗	✓	✓	✗	✗	✗
Public DGS Corpus [15]	DGS	–	50	327	✓	✓	✓	✓	✗	✗
BSL Corpus [31]	BSL	5k	–	249	✗	✓	✓	✗	✗	✗
Boston104 [39]	ASL	104	8.7 (min)	3	✗	✓	✓	✗	✗	✗
NCSLGR [24]	ASL	1.8k	5.3	4	✓	✓	✓	✗	✗	✗
How2Sign (ours)	ASL	16k	79	11	✓	✓	✓	✓	✓	✓

Table 1: **Summary of publicly available continuous sign language datasets.** To the best of our knowledge, How2Sign is the largest publicly available Sign Language dataset across languages in terms of vocabulary, as well as the largest American Sign Language (ASL) dataset in terms of video duration. We also see that How2Sign is the dataset with the most parallel modalities. A detailed explanation of each modality can be found in the subsection 3.2

controlled environments with a single RGB camera. Recent works in neural machine translation [8] and production [30, 28] have adopted *RWTH-Phoenix-2014T* [12], a dataset of German Sign Language (DGS) on the specific domain of weather forecast from a TV broadcast that features 9 signers. The *Public DGS Corpus* [15] and the *Video-Based CSL (Chinese Sign Language)*[16] provide much larger video collections enriched with the body keypoint of the signers. In the case of *Public DGS Corpus*, these are 2D poses estimated with OpenPose [9] and from different view points, while *Video-Based CSL* provides 3D joints and depth information thanks to the recordings with a Kinect camera. If we focus on American Sign Language (ASL), *RWTH-BOSTON-104* [39] only contains 8.7 minutes of grayscale video, while *NCSLGR* [24] is larger but an order of magnitude smaller than How2Sign. In terms of annotation, all datasets but *Video-Based CSL* provide gloss annotations, that is, a text-based transcription of the signs that can serve as a proxy in translation tasks.

Table 1 presents an overview of publicly available continuous sign language datasets ordered by vocabulary size³. An important factor for the lack of large-scale datasets is that the collection and annotation of continuous sign language data is a laborious and expensive task. It requires linguistic experts working together with a native speaker, e.g. a Deaf person. *RWTH-Phoenix-2014T* [12] is one of the few datasets that are publicly available and has been used for training deep neural networks. A recent re-alignment in the annotations also allows studying sign language translation. However, their videos cover just 11 hours of data from weather broadcasts, and are restricted to one domain.

In summary, the current publicly available datasets are constrained by one or more of the following: (i) limited vocabulary size, (ii) short video or total duration and (iii)

limited domain. The How2Sign dataset provides a considerably larger vocabulary than the existing ones, and it does so in the continuous sign language setting for a broader domain of discourse. It also is the first sign language dataset that contains speech thanks to its alignment with the existing How2 dataset [27].

3. The How2Sign dataset

The How2Sign dataset consists of a parallel corpus of speech and transcriptions of instructional videos and their corresponding American Sign Language (ASL) translation videos and annotations. A total of *80 hours* of multiview American Sign Language videos were collected, as well as gloss annotations [22] and a coarse video categorization.

Source language. The instructional videos translated into ASL come from the existing *How2 dataset* [27], a publicly available multimodal dataset for vision, speech and natural language understanding, with utterance-level time alignments between the speech and the ground-truth English transcription. Following the same splits from the *How2-300h* dataset, we selected a 60-hour subset from the training set and the complete validation and test sets to be recorded.

3.1. Sign language video recordings

Signers. In total, 11 people appear in the sign language videos of the How2Sign dataset; we refer to them as *signers*. Of the 11 signers, 5 self-identified as hearing, 4 as Deaf and 2 as hard-of-hearing. The signers that were hearing were either professional ASL interpreters (4) or ASL fluent.

Recording pipeline. The signer would first watch the video with the transcript as subtitles in order to become familiar with the overall content; this enables them to perform a richer translation. ASL translation videos were then recorded, while the signer was watching the video with subtitles, and at a slightly slower-than-normal (0.75) speed. For each hour of video recorded, the preparation, recording and

³An extended overview of related datasets can be found at: https://how2sign.github.io/related_datasets.html

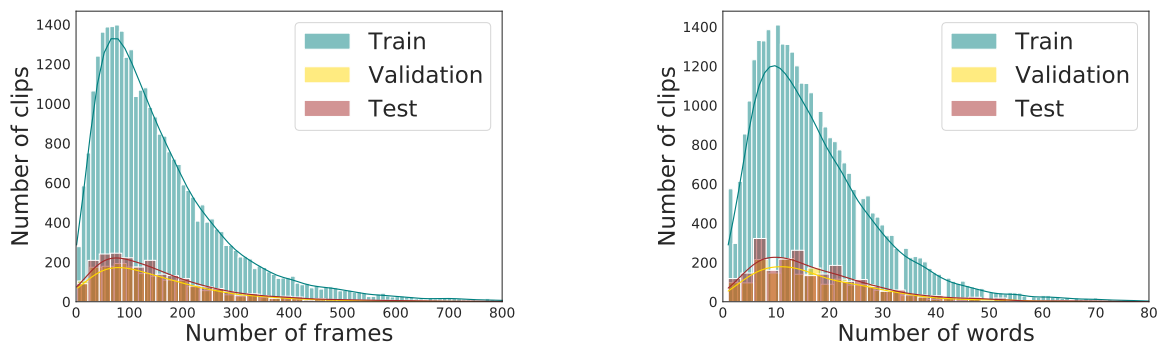


Figure 2: Distribution of the number of frames (left) and words (right) over sentence-level clips.

video review took approximately 3 hours on average.

All recordings were performed in a supervised setting in two different locations: at the *Green Screen studio* and at the *Panoptic studio*, both presented below. We recorded the complete 80 hours of the dataset in the green screen studio. We then chose a small subset of videos (approx. 3 hours) from the validation and test splits and recorded them again in the Panoptic studio. After recording, we trimmed all sign language videos and divided them in *sentence-level clips*, each annotated with a corresponding English transcript, and the modalities presented in Section 3.2.

Green screen studio. The *Green Screen studio* was equipped with a depth and a high definition (HD) camera placed in a frontal view of the participant, and another HD camera placed at a lateral view. All three cameras recorded videos at 1280x720 resolution, at 30 fps. Samples of data recorded in this studio are shown in the top row of Figure 1.

Panoptic studio. The *Panoptic studio* [17] is a system equipped with 480 VGA cameras, 30 HD cameras and 10 RGB-D sensors, all synchronized. All cameras are mounted over the surface of a geodesic dome⁴, providing redundancy for weak perceptual processes (such as pose detection) and robustness to occlusion. In addition to the multiview VGA and HD videos, the recording system can further estimate high quality 3D keypoints of the interpreters, also included in How2Sign. Samples of data recorded in this studio are shown on the bottom-right of Figure 1.

3.2. Dataset Modalities

The modalities enumerated in the columns of Table 1 are detailed in this section. Apart from the English translations and speech modalities that were already available from the How2 [27] dataset, all other modalities were either collected or automatically extracted. To the best of our knowledge, How2Sign is the largest publicly available sign

language dataset across languages in terms of vocabulary, as well as an order of magnitude larger than any other ASL dataset in terms of video duration. We see that How2Sign is also the dataset with the most parallel modalities, enabling multimodal learning.

Multiview. All 80 hours of sign language videos were recorded from multiple angles. This allows the signs to be visible from multiple points of view, reducing occlusion and ambiguity, especially in the hands. Specifically, the sign language videos recorded in the Green Screen studio contain two different points of view, while the Panoptic studio recordings consist of recordings of more than 500 cameras allowing for a high quality estimation of 3D keypoints [17].

Transcriptions. The English translation modality originates from the subtitles track of How2 original videos. The transcriptions were provided by the uploader of the instructional video in form of text, that was loosely synced with the video’s speech track. As subtitles are not necessarily fully aligned with the speech, transcriptions were time-aligned at the sentence-level as part of the How2 dataset [27].

Gloss is used in linguistics to transcribe signs using spoken language words. It is generally written in capital letters and indicates what individual parts of each sign mean, including annotations that account for facial and body grammar. An example of gloss annotation is shown on the bottom right of Figure 1. It is important to note that gloss is not a true translation, it instead provides the appropriate spoken language morphemes that express the meaning of the signs in spoken language [20, 22]. Glosses do not indicate special hand-shape, hand movement/orientation, nor information that would allow the reader to determine how the sign is made, or what its exact meaning in a given context. They also do not indicate grammatical uses of facial expressions (for example, raising the eyebrows is used in yes/no questions). Gloss is the form of text that is closest to sign language and it has been used by a number of approaches as an intermediate representation for sign language processing [12, 30, 28, 40, 19].

⁴<http://www.cs.cmu.edu/~hanbyulj/panoptic-studio/>

Pose information. Human pose information, *e.g.* body, hand and face keypoints were extracted for all the recorded sign language videos in the full resolution – 1280 x 720 pixels. For the Green Screen studio data, the 2-dimensional (2D) pose information was automatically extracted using OpenPose [9]. In total, each pose consists of 25 body keypoints, 70 facial keypoints and 21 keypoints for each hand. We provide pose information for both frontal and side view of the Green Screen studio data. A sample of the pose information extracted can be seen on the bottom row in the left side of Figure 1. For the Panoptic studio data, we provide high quality 3-dimensional (3D) pose information estimated by the Panoptic studio internal software [17] that can be used as ground-truth for a number of 3D vision tasks.

Depth data. For the Green Screen studio data, the sign language videos were also recorded using a Depth sensor (Creative BlasterX Senz3D) from the frontal viewpoint. The sensor has high precision facial and gesture recognition algorithms embedded and is able to focus on the hands and face, the most important human parts for sign language.

Speech. The speech track comes from the instructional videos as part of the How2 dataset [27].

3.3. Collected Annotations

Beyond the video recordings and automatically extracted pose information, we further collected a number of manual annotations for the sign language videos.

Gloss and sentence boundaries. We collected gloss annotations by employing ASL linguists. The annotations were collected using ELAN [13], an annotation software for audio and video recordings, specifically enhanced for sign language annotations. Information in ELAN is represented in tiers which are time-aligned to the video files, giving us the start and end boundaries of each sentence and producing what we call the sentence boundaries. The gloss annotation took in average one hour per 90 seconds of video.

Video Categories. Although the How2 dataset provides automatically extracted “topics” for all videos using Latent Dirichlet Allocation [5], we found that the automatic annotations were in general very noisy and not properly characterizing the selected videos. In order to better categorize the videos, we manually selected 10 categories⁵ from the instructional website Wikihow⁶ and manually classified each How2Sign video in a single category. The distribution of videos across the ten categories can be seen in Figure 3.

3.4. Dataset statistics

In Table 2 we show detailed statistics of the How2Sign dataset. A total of 2,456 videos from the How2 [27] were

⁵The categories are: Personal Care and Style, Games, Arts and Entertainment, Hobbies and Crafts, Cars and Other, Vehicles, Sports and Fitness, Education and Communication, Food and Drinks, Home and Garden and Pets and Animals.

⁶<https://www.wikihow.com/Special:CategoryListing>

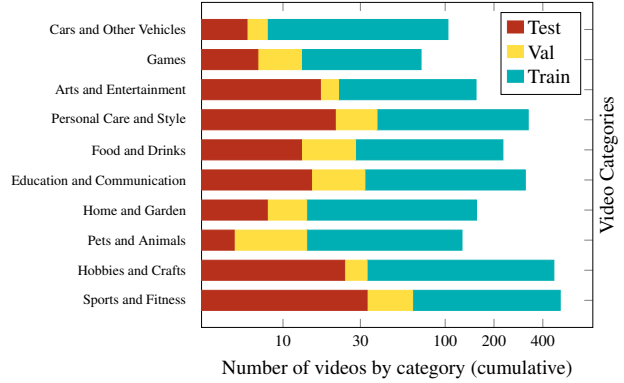


Figure 3: Cumulative number of videos per category.

used to record the sign language videos. Some of the videos were recorded more than once by a different signer in the Green screen studio – 21 videos from the training set, 17 videos from the validation set and 35 videos from the test set. All the recorded Videos were split into sentence-level clips. Each clip has on average 162 frames (5.4 seconds) and 17 words. The distribution of frames (right) and words (left) over all the clips for the 3 splits of the dataset can be seen in Figure 2. The collected corpus covers *more than 35k sentences* with an English vocabulary of more than 16k words. Where approximately, 20% of it is finger spelled. The videos were recorded by 11 different signers distributed across the splits. The test set contains 26 duplicated videos that were recorded by a signer that is not present in the training set; this subset of 26 videos can be used for measuring *generalization across different signers*. In total, 9 signers participated in the Green Screen studio recordings, and 6 signers in the Panoptic studio recordings. The bottom section of Table 2 refers to the automatically extracted human pose annotations.

3.5. Privacy, Bias and Ethical Considerations

In this section we discuss some metadata that we consider important for understanding the biases and generalization of the systems trained on our data.

Privacy. Since facial expressions are a crucial component for generating and/or translating Sign Language, it was not possible to avoid recordings that include the signer’s face. To that end, all the research steps followed procedures approved by the Carnegie Mellon University Institutional Review Board including a Social & Behavioral Research training done by the first and second authors, and a consent form provided by the participants agreeing on being recorded and making their data publicly available for research purposes. It is important to note that this puts at risk the authenticity of the linguistic data collected, as signers may monitor their production more carefully than usual.

Audiological status and language variety. The majority

	Green screen studio				Panoptic studio			
	train	val	test	Total	val	test	Total	
How2 [27] videos	2,192	115	149	2,456	48	76	124	
Sign language Videos	2,213	132	184	2,529	48	76	124	
Sign language video Duration (h)	69.62	3.91	5.59	79.12	1.14	1.82	2.96	
Number of frames (per view)	6.3M	362,319	521,219	7.2M	123,120	196,560	319,680	
Number of clips	31,128	1,741	2,322	35,191	642	940	1,582	
Camera views per SL video	3 HD + 1 RGB-D				480 VGA + 30 HD + 10 RGB-D			
Sentences	31,128	1,741	2,322	35,191	642	940	1,582	
Vocabulary size	15,686	3,218	3,670		1807	2360	3260	
Out-of-vocabulary	–	413	510					
Number of signers	8	5	6	9	3	5	6	
Signers not in train set	–	0	1		2	2		
		<u>2D keypoints</u>			<u>3D keypoints</u>			
Body pose		25				25		
Facial landmarks		70		137		70		
Hand pose (two hands)		21 + 21				21 + 21		

Table 2: Statistics of the **How2Sign** dataset. Some of the videos were recorded more than once by a different signer in the Green screen studio (see second row vs. first row). ASL videos recorded were split into sentence-level *clips*. Each clip has on average 162 frames (5.4 seconds) and 17 words.

of the participants identified American Sign Language and contact signing (Pidgin Sign English - PSE) as the main language used during the recordings. It is noteworthy that differences in audiological status are correlated with different language use. The Deaf were likely to identify ASL as the main language used in the recording process. In contrast, the hearing were likely to identify a mix of contact signing and ASL as the main language use in the recording process. More information about PSE and ASL can be found in [26].

Geographic. All participants were born and raised in the United States of America, and learned American Sign Language as their primary or second language at school time.

Signer variety. Our dataset was recorded by signers with different body proportions. Six of them were self-identified male and five self-identified female. The dataset was collected across 65 days during 6 months which gives a variety of clothing and accessories used by the participants.

Data bias. Our data does not contain large diversity in race/ethnicity, skin tone, background scenery, lighting conditions and camera quality.

4. Evaluating the potential of How2Sign for sign language tasks

The communication barrier between sign and non-sign language users may be reduced in the coming years thanks to the recent advances in neural machine translation and computer vision. Recent works are making steps towards sign language production [30, 33, 41, 40, 29] by automatically generating detailed human pose keypoints from spoken language, and translation [19], *i.e.*, using keypoints as

input to generate text.

While keypoints can carry detailed human pose information and can be an alternative for reducing the computational bottleneck that is introduced when working with the actual video frames, no studies have been made so far on whether they are indeed useful when it comes to understanding sign language by its users. In this section we present a study where we try to understand *if and how well sign language users understand automatically generated sign language videos* that use keypoints from How2Sign as sign language representation. We run this study with four ASL speakers and record their understanding of the generated videos in terms of the category, translation into American English, and a final subjective rating about how understandable the videos were.

4.1. Synthesizing sign language videos

We experiment with two ways of generating sign language videos: 1) skeleton visualizations and 2) Generative Adversarial Network generated (GAN-generated) videos.

Skeleton visualizations. Given a set of estimated keypoints, one can visualize them as a wired skeleton connecting the modeled joints (see the middle row of Figure 4).

GAN-generated videos. Another option would be to go one step further and use generative models to synthesize videos on top of predicted keypoints. To generate the animated video of a signer given a set of keypoints, we use the motion transfer and synthesis approach called Everybody Dance Now (EDN) [11]. This model is based on Pix2PixHD [36], but is further enhanced with a learned model of temporal coherence for better video and motion

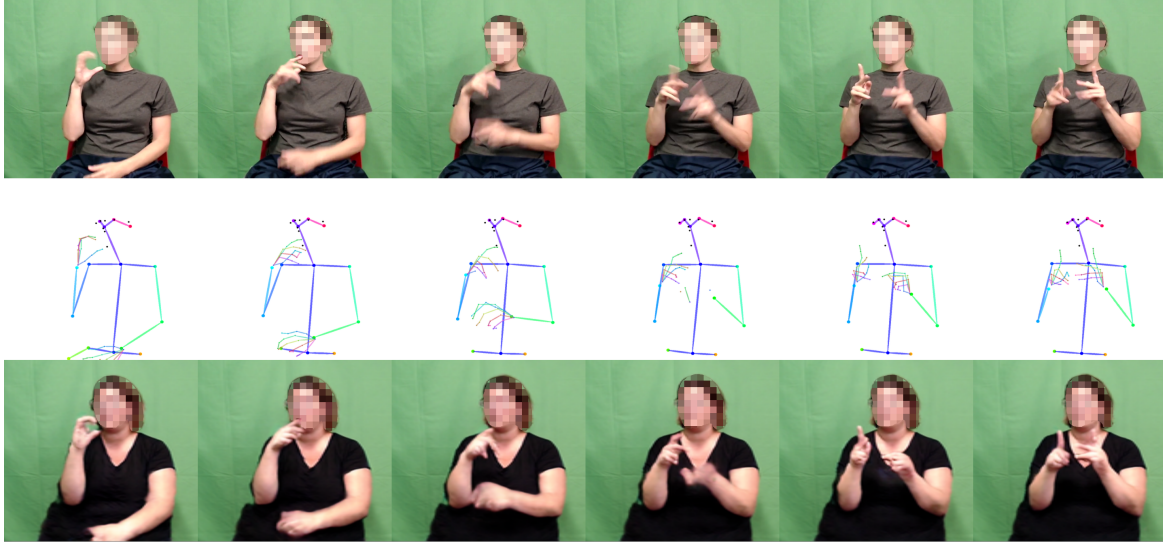


Figure 4: Sample of generated SL videos. Source video (top row) was used to automatically extract 2D keypoints (middle row) and generate frames of a video with a different identity (bottom row).

synthesis between adjacent frames by predicting two consecutive frames, as well as a separate module for high resolution face generation. It is worth noting that this approach models facial landmarks separately, something highly desirable in our case because they are one of the critical features for sign language understanding. The EDN model was trained on a subset of the How2Sign dataset that contains videos from *two* participants. Specifically, keypoints extracted from videos of the first signer (top row in Figure 4) were used to learn the model that generates realistic videos of the second signer (bottom row)⁷. The subset used consists of 28 hours of the training split.

4.1.1 Quantitative evaluation of the GAN-generated sign language videos.

An approximate but automatic way of measuring the visual quality of the generated videos is by comparing the keypoints that can be reliably detected by OpenPose in the source and generated videos. We focus only on the 125 upper body keypoints which are visible in the How2Sign videos, and discard those from the legs. We use two metrics: a) the Percentage of Detected Keypoints (PDK), which corresponds to the fraction of keypoints from the source frame which were detected in the synthesized frame, and b) the Percentage of Correct Keypoints (PCK) [38], which labels each detected keypoint as “correct” if the distance to the keypoint in the original image is less than 20% of the torso diameter in all keypoints and 10% of the torso diameter for the hands.

⁷A sample of a generated video can be seen at: <https://youtu.be/wOxWUyXX6Ys>

OP confidence scores	PDK			PCK		
	0	0.2	0.5	0	0.2	0.5
All keypoints	0.99	0.88	0.87	0.90	0.94	0.96
Hands	0.99	0.38	0.17	0.08	0.11	0.12

Table 3: Percentage of Detected Keypoints (PDK) and Percentage of Correct Keypoints (PCK) for all keypoints and just for the hands, when thresholding at different detection confidence scores of OpenPose (OP).

In Table 3 we present these metrics for different minimum confidence thresholds of the OpenPose (OP keypoint detectors). We report results for all keypoints, as well as when restricting the evaluation only on the hand keypoints. We see that although the repeatability of keypoints is high in general, the model fails to predict reliable keypoints for the hands. This limitation is especially relevant in sign language processing.

4.2. Can ASL signers understand generated sign language videos?

We evaluate the degree of understanding for both skeleton visualizations and the GAN-generated videos by showing 3-minute-long videos to four ASL signers. Two of them watched the skeletons visualizations, while the other two watched the GAN-generated videos. During the evaluation, each subject was asked to: a) classify six videos between the ten video categories (see subsection 3.2 for more information about the dataset categories); b) answer the question “How well could you understand the video?” on the five-level scale ((1) Bad, (2) Poor, (3) Fair, (4) Good, (5) Excel-

	Acc.	MOS	BLEU-1	BLEU-2	BLEU-3	BLEU-4
Skeleton	83.3 %	2.50	10.90	3.02	1.87	1.25
GAN-generated	91.6 %	2.58	12.38	6.71	3.32	1.89

Table 4: Comparison between generated skeletons and GAN videos in terms of classification (Accuracy), mean opinion score (MOS) and translation (BLEU) [25].

GT	I'm not going to use a lot, I'm going to use very very little.
Skeleton	That is not too much don't use much, use a little bit
EDN	Don't use a lot, use a little dont use lot use little bit
GT	I'm going to dice a little bit of peppers here.
Skeleton	cooking chop yellow peppers
EDN	cook with a little pepper chop it little bit and sprinkle

Table 5: Ground-truth (GT) and collected translations for two clips of the ‘‘Food and Drink’’ category. All subjects were able to correctly classify the category.

lent); and c) watch two clips from the previously seen video and translate them into American English. Results averaged over all subjects are presented in Table 4. We report accuracy for the classification task, the Mean Opinion Score (MOS) for the five-scale question answers and BLEU [25] scores for the American English translations. Qualitative results are shown in Table 5.

Results show a preference towards the generated videos rather than the skeleton ones, as the former result in higher scores across all metrics. In terms of general understanding of the topic, the subjects were able to mostly classify the videos correctly with both types of visualizations.

When it comes to finer grained understanding measured via the English translations, however, we can see from both Table 4 and Table 5 that neither skeletons nor GAN-generated videos are sufficient to convey important information needed from ASL signers to completely understand the sign language sentences. We hypothesize that current human pose estimation methods such as [9] are still not mature enough when it comes to estimate fast movements of the hands. We observed that due to the nature of sign language and the fast movements of the signers’ hands, OpenPose lacks precision in those cases which can make the visualizations incomplete, harming the understanding of some important parts of sign language.

How can computer vision do better? Our results show that the EDN model used as an out-of-the-box approach is not enough for sign language video generation. Specifically, we show that the model struggles with generating the hands and detailed facial expressions, which play a central role in sign language understanding. We argue that human pose es-

timation plays an important key in this aspect and needs to be more robust to blurry images, especially in the hands and to fast movements in order to be suitable to sign language research. We also argue that it is worth pursuing generative models that focus on generating hand details, particularly on the movements of the fingers, as well as clear facial expressions on full-body synthesis.

5. Conclusion

In this paper, we present the How2Sign dataset, a *large-scale multimodal and multiview dataset of American Sign Language*. With more than 80 hours of sign language videos and their corresponding speech signal, English transcripts and annotations, How2Sign has the potential to impact a wide range of sign language understanding tasks, such as sign language recognition, translation and production, as well as wider multimodal and computer vision tasks like 3D human pose estimation. How2Sign extends the How2 [27] dataset, an existing multimodal dataset with a new sign language modality, and therefore enables connecting with research performed in the vision, speech and language communities. In addition to that, we further conducted a study in which sign language videos generated from the automatically extracted annotations of our dataset were presented to ASL signers. To our knowledge, this is the first study how well keypoint-based synthetic videos, a commonly used representation of sign language production and translation, can be understood by sign language users. Our study indicates that current video synthesis methods allow the understanding to a certain extent *i.e.*, the classification of the video category, but lack in fidelity to allow for a fine-grained understanding of the complete sign language sentence.

Acknowledgments

This work received funding from Facebook through gifts to CMU and UPC; through projects TEC2016-75976-R, TIN2015-65316-P, SEV-2015-0493 and PID2019-107255GB-C22 of the Spanish Government and 2017-SGR-1414 of Generalitat de Catalunya. This work used XSEDE’s ‘‘Bridges’’ system at the Pittsburgh Supercomputing Center (NSF award ACI-1445606). Amanda Duarte has received support from la Caixa Foundation (ID 100010434) under the fellowship code LCF/BQ/IN18/11660029. Shruti Palaskar was supported by the Facebook Fellowship program. The authors would like to thank Chinmay Hejmadi, Xabier Garcia and Brandon Taylor for their help during the data collection and processing and Yannis Kalantidis for his valuable feedback. This work would not be possible without the collaboration and feedback from the signers and the Deaf community involved throughout the project.

References

- [1] World Health Organization 2019. Deafness and hearing loss. <https://www.who.int/news-room/fact-sheets/detail/deafness-and-hearing-loss>. Accessed: 2019-05-19. 1
- [2] Samuel Albanie, Gül Varol, Liliane Momeni, Triantafyllos Afouras, Joon Son Chung, Neil Fox, and Andrew Zisserman. Bsl-1k: Scaling up co-articulated sign language recognition using mouthing cues. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [3] David F Armstrong, William C Stokoe, and Sherman E Wilcox. *Gesture and the nature of language*. Cambridge University Press, 1995. 1, 2
- [4] Vassilis Athitsos, Carol Neidle, Stan Sclaroff, Joan Nash, Alexandra Stefan, Quan Yuan, and Ashwin Thangali. The american sign language lexicon video dataset. In *CVPRW'08.*, pages 1–8. IEEE, 2008. 2
- [5] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003. 5
- [6] Danielle Bragg, Oscar Koller, Mary Bellard, Larwan Berke, Patrick Boudreault, Annelies Braffort, Naomi Caselli, Matt Huenerfauth, Hernisa Kacorri, Tessa Verhoef, et al. Sign language recognition, generation, and translation: An interdisciplinary perspective. In *The 21st International ACM SIGACCESS Conference on Computers and Accessibility*, pages 16–31, 2019. 1, 2
- [7] Ruth Butler, Tracey Skelton, and Gill Valentine. Language barriers: Exploring the worlds of the deaf. *Disability Studies Quarterly*, 21(4), 2001. 1
- [8] Necati Cihan Camgoz, Oscar Koller, Simon Hadfield, and Richard Bowden. Sign language transformers: Joint end-to-end sign language recognition and translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10023–10033, 2020. 3
- [9] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *arXiv preprint arXiv:1812.08008*, 2018. 3, 5, 8
- [10] N Caselli, Z Sevcikova, A Cohen-Goldberg, and K Emmorey. Asl-lex: A lexical database for asl. *Behavior Research Methods*, 2016. 2
- [11] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei Efros. Everybody dance now. In *ICCV*, 2019. 6
- [12] Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. Neural sign language translation. In *CVPR*, pages 7784–7793, 2018. 2, 3, 4
- [13] Onno Crasborn and Sloetjes Han. Enhanced elan functionality for sign language corpora. *Journal of deaf studies and deaf education*, 2008. 5
- [14] Warren R Goldmann and James R Mallory. Overcoming communication barriers: communicating with deaf people. 1992. 1
- [15] Thomas Hanke, Marc Schuder, Reiner Konrad, and Elena Jahn. Extending the public dgs corpus in size and depth. In *LREC2020 - Workshop on the Representation and Processing of Sign Languages*, pages 75–82, 2020. 3
- [16] Jie Huang, Wengang Zhou, Qilin Zhang, Houqiang Li, and Weiping Li. Video-based sign language recognition without temporal segmentation. In *AAAI*, 2018. 3
- [17] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3334–3342, 2015. 2, 4, 5
- [18] Hamid Reza Vaezi Joze and Oscar Koller. Ms-asl: A large-scale data set and benchmark for understanding american sign language. *arXiv preprint arXiv:1812.01053*, 2018. 2
- [19] Sang-Ki Ko, Chang Jo Kim, Hyedong Jung, and Choongsang Cho. Neural sign language translation based on human key-point estimation. *Applied Sciences*, 9(13), 2019. 2, 4, 6
- [20] Jolanta Lapiak. Gloss: transcription symbols. <https://www.handspeak.com/learn/index.php?id=3>. Accessed: 2019-08-20. 4
- [21] Dongxu Li, Cristian Rodriguez, Xin Yu, and Hongdong Li. Word-level deep sign language recognition from video: A new large-scale dataset and methods comparison. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1459–1469, 2020. 2
- [22] Scott K Liddell et al. *Grammar, gesture, and meaning in American Sign Language*. Cambridge University Press, 2003. 2, 3, 4
- [23] Aleix M Martínez, Ronnie B Wilbur, Robin Shay, and Avinash C Kak. Purdue rvl-slll asl database for automatic recognition of american sign language. In *Proceedings. Fourth IEEE International Conference on Multimodal Interfaces*, pages 167–172. IEEE, 2002. 2
- [24] Carol Neidle and Christian Vogler. A new web interface to facilitate access to corpora: Development of the asllrp data access interface (dai). In *Proc. 5th Workshop on the Representation and Processing of Sign Languages: Interactions between Corpus and Lexicon, LREC*, 2012. 3
- [25] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic evaluation of machine translation. In *ACL*, 2002. 8
- [26] Judy Reilly and Marina L McIntire. American sign language and pidgin sign english: What’s the difference? *Sign Language Studies*, pages 151–192, 1980. 6
- [27] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze. How2: a large-scale dataset for multimodal language understanding. *arXiv preprint arXiv:1811.00347*, 2018. 2, 3, 4, 5, 6, 8
- [28] Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. Adversarial training for multi-channel sign language production. In *The 31st British Machine Vision Virtual Conference (BMVC)*, 2020. 3, 4
- [29] Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. Everybody sign now: Translating spoken language to photo realistic sign language video. *arXiv preprint arXiv:2011.09846*, 2020. 6
- [30] Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. Progressive transformers for end-to-end sign language

- production. In *European Conference on Computer Vision (ECCV)*, 2020. [2](#), [3](#), [4](#), [6](#)
- [31] Adam Schembri, Jordan Fenlon, Ramas Rentelis, Sally Reynolds, and Kearsy Cormier. Building the british sign language corpus. *Language Documentation & Conservation*, 7:136–154, 2013. [2](#), [3](#)
- [32] William C Stokoe Jr. Sign language structure: An outline of the visual communication systems of the american deaf. *Journal of deaf studies and deaf education*, 10(1):3–37, 2005. [2](#), [11](#)
- [33] Stephanie Stoll, Necati Cihan Camgoz, Simon Hadfield, and Richard Bowden. Text2sign: towards sign language production using neural machine translation and generative adversarial networks. In *International Journal of Computer Vision*, 2020. [2](#), [6](#)
- [34] Špela Vintar, Boštjan Jerko, and Marjetka Kulovec. Compiling the slovene sign language corpus. In *5th Workshop on the Representation and Processing of Sign Languages: Interactions between Corpus and Lexicon. Language Resources and Evaluation Conference (LREC)*, volume 5, pages 159–162, 2012. [2](#)
- [35] U. Von Agris and K.-F. Kraiss. Signum database: Video corpus for signer-independent continuous sign language recognition. In *Workshop on Representation and Processing of Sign Languages*, pages 243–246, 2010. [2](#), [3](#)
- [36] 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*, 2018. [6](#)
- [37] James C Woodward. Implications for sociolinguistic research among the deaf. *Sign Language Studies*, pages 1–7, 1972. [1](#)
- [38] Yi Yang and Deva Ramanan. Articulated human detection with flexible mixtures of parts. *IEEE TPAMI*, 35:2878–90, 12 2013. [7](#)
- [39] Morteza Zahedi, Philippe Dreuw, David Rybach, Thomas Deselaers, and Hermann Ney. Continuous sign language recognition—approaches from speech recognition and available data resources. In *Workshop on Representation and Processing of Sign Languages*, 2006. [3](#)
- [40] Jan Zelinka and Jakub Kanis. Neural sign language synthesis: Words are our glosses. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 3395–3403, 2020. [2](#), [4](#), [6](#)
- [41] Jan Zelinka, Jakub Kanis, and Petr Salajka. Nn-based czech sign language synthesis. In *International Conf. on Speech and Computer*, pages 559–568. Springer, 2019. [2](#), [6](#)

Supplementary Material

6. Sign Language

In this section we discuss in more detail some important non-manual features (that are not conveyed through other linguistic parameters *e.g.* palm orientation, handshape, etc.) present in sign languages. It is important to remember that American Sign Language, for example, requires more than just complex hand movements to convey a message. Without the use of proper facial expressions and other non-manual features as the ones described below, a message could be greatly misunderstood [32].

Head movement. The movement of the head supports the semantics of sign language. Questions, affirmations, denials, and conditional clauses are communicated with the help of the signer’s head movement.

Facial grammar. Facial grammar does not only reflect a person’s affect and emotions, but also constitutes to large part of the grammar in sign languages. For example, a change of head pose combined with the lifting of the eye brows corresponds to a subjunctive.

Mouth morphemes (mouthing). Mouth movement or mouthing is used to convey an adjective, adverb, or another descriptive meaning in association with an ASL word. Some ASL signs have a permanent mouth morpheme as part of their production. For example, the ASL word NOT-YET requires a mouth morpheme (TH) whereas LATE has no mouth morpheme. These two are the same sign but with a different non-manual signal. These mouth morphemes are used in some contexts with some ASL signs, not all of them.

7. How2Sign dataset

Here we discuss some additional metadata that are important for a better understanding of our data as well as the biases and generalization of the systems trained using the How2Sign dataset. We also describe information that might be helpful for future similar data collection.

Gloss. We collected gloss annotations for the ASL videos present in the How2Sign dataset using ELAN. Figure S2 shows samples of the gloss annotations present in our dataset. Here we describe some conventional and few modified symbols and explanations that will be found in our dataset. A complete list is available on the dataset website.

- *Capital letters.* English glosses are written using capital letters. They represent an ASL word or sign. It is important to remember that gloss is not a translation. It is only an approximate representation of the ASL sign itself, not necessarily a meaning.
- A *hyphen* is used to represent a single sign when more than one English word is used in gloss (*e.g.* STARE-AT).
- The *plus sign (+)* is used in ASL compound words (*e.g.* MOTHER+FATHER – used to transcribe parents). It is also used when someone combines two signs in one (*e.g.* YOU THERE will be glossed as YOU+THERE).
- The *plus sign (++)* at the end of a gloss indicates a number of repetitions of an ASL sign (*e.g.* AGAIN++ – the word “again” was signed two more times meaning “again and again”).
- *FS:* represents a fingerspelled word (*e.g.* FS:AMELIA).

- *IX* is a shortcut for “index”, which means to point to a certain location, object, or person.
- *LOC* is a shortcut for “locative”, a part of the grammatical structure in ASL.
- *CL:* is a shortcut for “classifier”. Classifiers are signs that use handshapes that are associated with specific categories (classes) of things, size, shape, or usage. They can help to clarify the message, highlight specific details, and provide an efficient way of conveying information⁸. In our annotations, classifiers will appear as: “CL:classifier(information)”. For example, if the signer signs “TODAY BIKE” and uses a classifier to show the bike going up the hill, this would be glossed as: “TODAY BIKE CL:3 (going uphill)”.

Signers. Figure S1 show all the 11 signers that participated in the recordings of the How2Sign dataset. From the 11 signers, four of them (signers 1, 2, 3 and 10) participated in both the Green Screen studio and the Panoptic studio recordings. Signers 6 and 7 participated only in the Panoptic studio recordings, while signers 4, 5, 8, 9 and 11 participated only in the Green Screen recordings. The signer ID information of each video is also made available.

Recording pipeline. *Importance of providing the speech and original video to the signer before the recordings:* As part of the design phase of our data collection, signers were asked to perform English to ASL translation when given: (1) just text without reading it beforehand; (2) the video and text together but without seeing it previously and (3) text and video together and allowing them to watch it before the recording. The conclusions for each case were: (1) signers found it hard to understand and follow the lines at the same time, causing lots of pauses and confusion; (2) signers found it easier to understand and translate but still with some pauses and (3) the understanding and flow improved.

7.1. Discussion

How high is the quality of the extracted keypoints? We conducted a number of studies to estimate the quality of the automatically extracted 2D poses. A number of sanity checks showed us that extracting keypoints in higher resolution (1280 x 720) resulted to pose estimation that have on average higher confidence – 53.4% average keypoint confidence for high resolution versus 42.4% confidence for low resolution (210 x 260). This difference is more prominent when different parts of the body are analyzed. Table S1 show the different average confidence scores when OpenPose is extracted using high and low resolution videos. We see that both hands are the most harm when low resolution is used.

More importantly, in Section 4 we present a study with native speakers and verified that our 2D keypoints are sufficient to a certain degree for sign language users to classify and transcribe the ASL videos back to English.

Factors that may impair accurate automatic tracking. During the recording, signers were requested to not use loose clothes, rings, earrings, watch, or any other accessories that might impair accurate automatic tracking. They were also asked to wear solid colored shirts (that contrast with their skin tone).

Out-of-vocabulary and signer generalization. Although not

⁸More info about handshapes and classifiers can be found at: <https://www.lifefprint.com/asl101/pages-signs/classifiers/classifiers-main.htm>



Figure S1: All the 11 signers that appear in the How2Sign dataset videos. On the top row, we can see signers 1-5 (from left to right) in the Green Screen Studio, while on the bottom row we can see signers 8-11 (again left to right) in the Green Screen Studio. The rightmost figure on the bottom row shows signers 6-7 in the Panoptic studio.

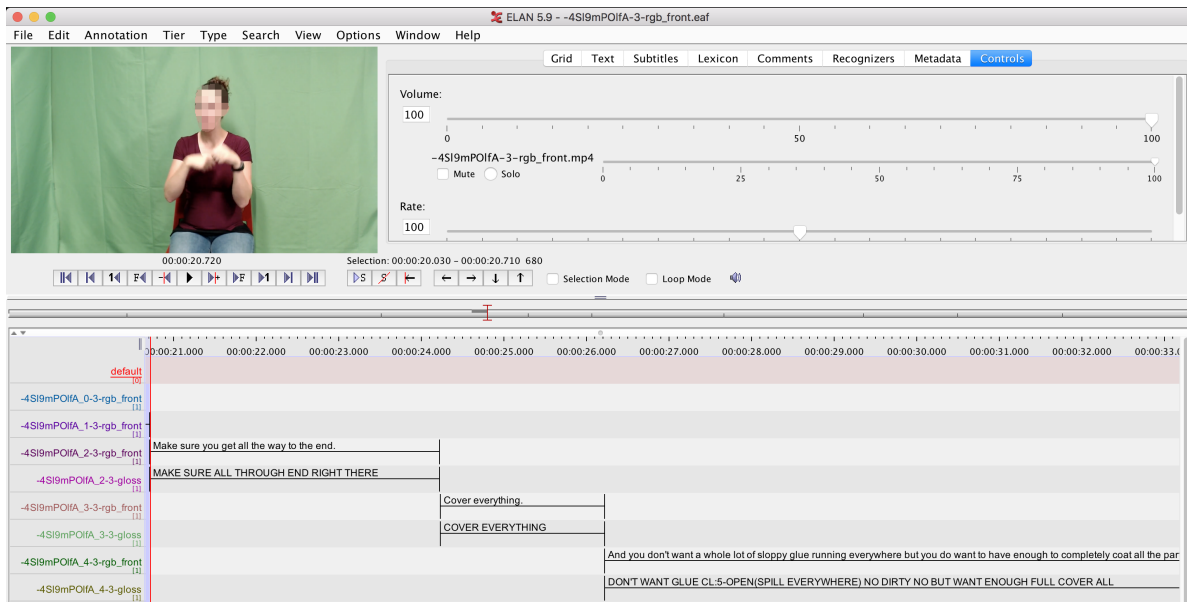


Figure S2: Samples of gloss annotations collected using ELAN.

	Body	Right hand	Left hand	Face	Total
High resolution	0.39	0.42	0.47	0.84	0.53
Low resolution	0.40	0.24	0.30	0.73	0.42

Table S1: Average of confidence score of OpenPose on high resolution (1280 x 720) compared with low resolution (210 x 260) videos of the How2Sign dataset.

specifically designed for this, the How2Sign dataset can be used for measuring generalization with respect to both out-of-vocabulary words and signers. The dataset contains 413 and 510 out-of-vocabulary words, *e.g.* words that occur in validation and test, respectively, but not in training. It further contains duplicate

recordings on the test set by a signer that is *not present in the training set*; these recordings can be used for measuring generalization across different signers and help understand how well the models can recognise or translate the signs given an out of the distribution subject.

Language variety. As discussed in subsection 3.5 our dataset contains variations in the language used during the recordings by each signer. In addition to that, we also would like to mention that sign language speakers can also use different signs or different linguistic registers (*i.e.*, formal or casual) to express the same given sentence. As we can see in Figure S3, two signers from our dataset used two different signs in a linguistic register to express the phrase “I am”. The signer on the left used the casual approach of signing (ME NAME) while the signer on the left used the formal approach (ME).

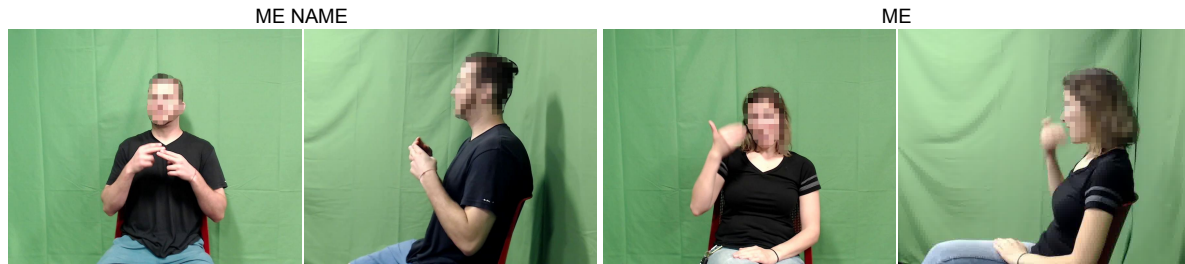


Figure S3: Sample of language variety on our dataset. Both signers were translating the sentence “I am”. We can see that the signer on the left used the casual approach of signing it (ME NAME) while the signer on the right used the formal approach (ME).

Intra-sign variety. In addition to the variety of signs and linguistic registers, it is also common to notice differences in the way of performing the same sign. For example, we can see on Figure S4 two signers from our dataset signing the word “hair”. In this sign, as described by its gloss annotation (IX-LOC-HAIR) the signer points to their own hair location. While performing the sign, the person can use slightly different locations to point at.

7.2. How2Sign statistics per signer

Table S2 presents detailed statistics of the videos from the How2Sign dataset recorded in the *Green Screen studio* grouped by signer.



Figure S4: Sample of intra-sign variety. In this case, both signers are signing the word “hair” (IX-LOC-HAIR). We can see that the on the left choose to point to her hair on a different position from the signer on the right.

	Signer 1	Signer 2	Signer 3	Signer 4	Signer 5	Signer 8	Signer 9	Signer 10	Signer 11	Total
Train										
Videos	50	22	163	24	899	994	18	-	43	2213
Hours	1.89	0.82	3.80	0.82	31.59	28.28	0.67	-	1.72	69.59
Utterances	892	422	1859	398	12102	14596	292	-	486	31047
Test										
Videos	16	16	37	-	47	42	-	26	-	184
Hours	0.51	0.53	1.05	-	1.67	1.08	-	0.71	-	5.55
Utterances	224	243	538	-	621	449	-	268	-	2343
Validation										
Videos	17	19	27	-	37	32	-	-	-	132
Hours	0.57	0.68	0.65	-	1.20	0.79	-	-	-	3.89
Utterances	276	270	306	-	454	433	-	-	-	1739

Table S2: Statistics of the *Green Screen studio* data by signer. We present the number of videos recorded by signer (videos), together with the total duration of the recorded videos in hours (Hours) and the number of utterances (Utterances) of each signer.