

Speak, Decipher and Sign: Toward Unsupervised Speech-to-Sign Language Recognition

Anonymous ACL submission

Abstract

Existing supervised sign language recognition systems rely on an abundance of well-annotated data. Instead, an unsupervised speech-to-sign language recognition (SSR-U) system learns to translate between spoken and sign languages by observing only non-parallel speech and sign-language corpora. We propose speech2sign-U, a neural network-based approach capable of both character-level and word-level SSR-U. Our approach significantly outperforms baselines directly adapted from unsupervised speech recognition (ASR-U) models by as much as 50% recall@10 on several challenging American sign language corpora with various levels of sample sizes, vocabulary sizes, and audio and visual variability. The code is available [here](#).

1 Introduction

Many hearing-impaired people communicate natively in sign language (SL); for them, SL communication is as effortless as native spoken communication is for normal-hearing people. However, when it comes to a conversation between a hearing-impaired and a normal hearing, tremendous barriers exist for several reasons. First, there is a shortage of people who are bilingual in spoken and sign languages. Automatic sign language recognition models exists (Koller et al., 2016; Huang et al., 2018) but are fully supervised and require a large number of annotated data, which are hard to acquire. As a result, such systems are often limited to a small vocabulary. On the other hand, untranscribed speech audio and SL videos are quite common on the Internet, presenting an exciting possibility: Given a non-parallel pair of speech and sign language datasets, can we train a model to translate between spoken and sign languages? This task, we called *unsupervised speech-to-sign language recognition* (SSR-U), is analogous to well-known problems such as unsupervised machine translation (MT-U) (Ravi and Knight, 2011; Artetxe et al.,

2018a; Lample et al., 2018) and unsupervised automatic speech recognition (ASR-U) (Liu et al., 2018; Chen et al., 2019; Baevski et al., 2021), albeit with a few new challenges. First of all, in the case of SSR-U, both modalities are continuous as opposed to at least one of them being discrete in the case of ASR-U and MT-U. Consequently, the matching process is much more challenging due to higher within and cross-modal variability. Further, most sign language and spoken language can only be matched on the *word* level as opposed to the sub-word level in the case of ASR-U. Not only does the space of possible mappings explode combinatorially, but less training data and fewer temporal constraints are also available to recover the correct mapping.

In this paper, we develop a neural network-based framework, speech2sign-U, for both character-level (with fingerspelling sequence) and word-level SSR-U. It achieves promising results on datasets with up to around 900 ASL signs.

2 Problem formulation

Suppose we have a corpus of unlabeled speech recordings sampled from the random process $A = (A_1, \dots, A_T)$ and another separately collected corpus of unlabeled sign language videos sampled from the random process $V = (V_1, \dots, V_L)$. Both A and V contain the *same* semantic information but different para-linguistic information such as speaker/signer identity and prosody. In other words, if we filter out the para-linguistic information and retain the semantic information as $X := X(A) = (X_1, \dots, X_T) \sim P_X$ for the speech and $Y := Y(V) = (Y_1, \dots, Y_L) \sim P_Y$ for the videos, we can find a *generator* function $\mathcal{G} : \mathbb{X}^T \mapsto \mathbb{Y}^L$ such that $Y = \mathcal{G}(X)$. Since the corpora are unpaired, we cannot estimate \mathcal{G} directly from samples, and the goal of SSR-U is to “decipher” it using only the relations between the speech-only and video-only distributions, P_X and

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017

018

019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

082 P_Y :

083
$$\sum_{x \in \mathbb{X}^T} P_X(x) G(y|x) = P_Y(y), \quad (1)$$

084 for all sign language unit sequences $y \in \mathbb{Y}^L$, where
085 $G(y|x) = 1$ if and only if $y = \mathcal{G}(x)$.086

3 Proposed Methods

087

3.1 Character-level speech2sign-U

088 In the case of character-level speech2sign-U, V is
089 drawn from a collection of unlabeled fingerspelling
090 sequences, where each V_i is the hand gesture for a
091 character. In this case, we adopt a similar architecture
092 as wav2vec-U (Baevski et al., 2021).093 **Sign video preprocessing** Given a sign video
094 $v \sim V$, we obtain its visual features (v_1, \dots, v_L)
095 by passing the raw video frames into a *local* feature
096 extractor such as VGG19 or RCNN (Ren et al.,
097 2015). The local features are then contextualized
098 by a *sign language encoder*, consisting of a two-
099 layer multilayer perceptron (MLP) and a one-layer
100 uni-directional LSTM:

101 $c_1, \dots, c_L = \text{LSTM}(v_1, \dots, v_L). \quad (2)$

102 The sign language encoder is then trained using
103 contrastive predictive coding (CPC) (van den Oord
104 et al., 2018):

105
$$\mathcal{L}_{\text{CPC}} := -\mathbb{E}_V \left[\sum_{i,k} \log \frac{e^{c_i^\top \text{MLP}(v_{i+k})}}{\sum_{n \in \mathcal{N}_{i,k}} e^{c_i^\top \text{MLP}(n)}} \right], \quad (3)$$

106 where $\mathcal{N}_{i,k}$ is a set of negative samples chosen uniformly at random from times other than $i + k$. Finally, we apply K-means clustering on (c_1, \dots, c_L) to obtain the *sign cluster units* $Y := (y_1, \dots, y_L)$.107 **Speech preprocessing** As in wav2vec-U, for
108 each utterance, we first use a voice activity de-
109 tector (VAD) to remove silences between speech
110 frames and randomly insert silences between word
111 boundaries of the sign cluster sequence so that their
112 silence distributions match. Next, we contextualize
113 the raw speech frames using wav2vec 2.0 pre-
114 trained on LibriLight:

115 $(z_1, \dots, z_T) = \text{wav2vec2}(a_1, \dots, a_T). \quad (4)$

116 Finally, we extract K-means clusters from
117 (z_1, \dots, z_T) and merge consecutive frames belong-
118 ing to the same clusters to obtain the segment-level
119 speech features (x_1, \dots, x_T) .120 **Unsupervised training** A convolutional gener-
121 ator $G : \mathbb{X} \rightarrow \mathbb{Y}$ then generates a sequence of
122 cluster units $(\hat{Y}_1, \dots, \hat{Y}_L) = \mathcal{G}(X)$ from the seg-
123 ment features X by sampling from the posterior
124 probabilities at each segment i :

125
$$\hat{Y}_i \sim G_i(y_i|X) := \frac{\exp(\text{Conv}_{i,y_i}(X))}{\sum_k \exp(\text{Conv}_{i,k}(X))}. \quad (5)$$

126 Then we adopt the generative adversarial network
127 (GAN (Goodfellow et al., 2014)) objective by train-
128 ing a binary classifier $D : \mathbb{Y} \mapsto [0, 1]$ to discrim-
129 inate between the real cluster sequence and the
130 generated one:

131
$$\begin{aligned} & \min_D \max_G -\mathbb{E}_{X \sim P_X} [\log(1 - D(G(X)))] \\ & \quad - \mathbb{E}_{Y \sim P_Y} [\log D(Y)] \\ & \quad + \lambda \mathcal{L}_{\text{gp}} + \gamma \mathcal{L}_{\text{sp}} + \eta \mathcal{L}_{\text{cd}}, \end{aligned} \quad (6)$$

132 where \mathcal{L}_{gp} , \mathcal{L}_{sp} and \mathcal{L}_{cd} stand for the gradient
133 penalty, smoothness penalty and code diversity
134 losses as defined in (Baevski et al., 2021).135

3.2 Word-level speech2sign-U

136 Word-level speech2sign-U is more challenging than
137 character-level: the GAN objective in Eq. (6) fails
138 to converge for vocabulary sizes of 100 or larger,
139 apparently due to variability in the audio and video
140 signals. Therefore, we instead adopt a novel GAN-
141 free architecture trained to match marginals be-
142 tween the generated and real probability distribu-
143 tions as shown in Fig. 1.144 **Preprocessing** We extract the sign video and
145 speech features similar to Section 3.1, except with a
146 few modifications: first, we assume the word-level
147 boundaries for both the speech and sign videos are
148 available, which may be ground truth or bound-
149aries detected using unsupervised word segmenta-
150 tion algorithms from phoneme boundaries (Kreuk
151 et al., 2020; Bhati et al., 2021; Cuervo et al., 2022).
152 Then we compute the segment-level speech fea-
153 tures by averaging the frame-level wav2vec 2.0
154 features within each word. Further, we use the
155 I3D (Carreira and Zisserman, 2017) as the local
156 feature extractor and average the pretrained video
157 feature frames within each word-level sign video
158 segment. Lastly, we perform K-means clustering
159 on the segment features and use the output cluster
160 units as inputs X to the speech generator as we
161 found that quantized speech features work better
162 than continuous features.

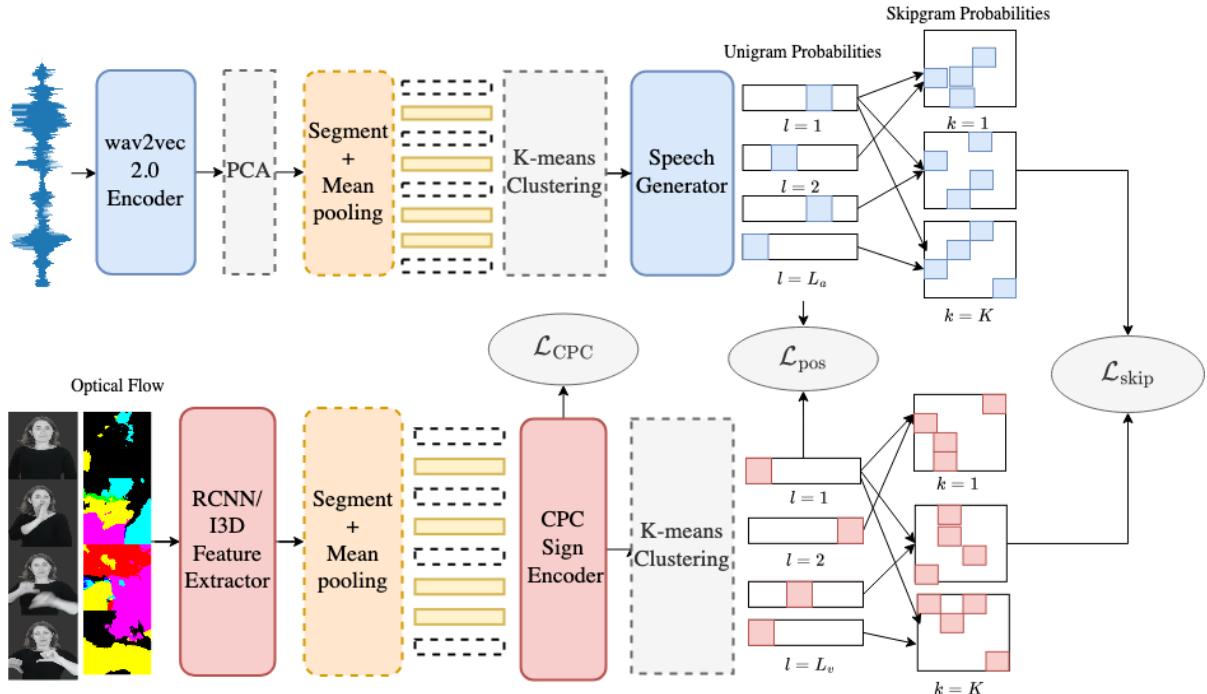


Figure 1: Overall architecture of the word-level speech2sign-U. Solid blocks contain trainable parameters while dashed blocks do not.

Unsupervised unigram matching Similar to Section 3.1, we seek to match the probability distributions in the two modalities as our unsupervised training criterion. Instead of using a convolutional generator as in Eq. (5), we instead use a *linear* generator for each segment i :

$$G(y_i|x_i) := \frac{\exp(W_{y_i}x_i)}{\sum_{y' \in \mathbb{Y}} \exp(W_{y'}x_i)}, \quad (7)$$

Eq. (1) can now be achieved by minimizing the ℓ_1 distance between the empirical *positional unigram* probabilities of the generated and real sign cluster units:

$$\mathcal{L}_{\text{pos}}(G) = \sum_{i=1}^L \|\hat{P}_{X_i}G - \hat{P}_{Y_i}\|_1, \quad (8)$$

where \hat{P}_{X_i} and \hat{P}_{Y_i} are empirical unigram distributions for the speech and sign units, and $G \in \mathbb{R}^{|\mathbb{X}| \times |\mathbb{Y}|} := (G(y|x))_{x \in \mathbb{X}, y \in \mathbb{Y}}$. Note that such an objective is typically optimized implicitly by a GAN, but we found that the explicit formula not only avoids the need for a discriminator but also leads to more stable training and better performance.

Unsupervised skip-gram matching Positional unigram constraints alone may not be sufficient for

word-level SSR-U. Therefore, we add additional moment constraints using *skip-grams*. Define the k -step skip-gram to be the joint probability

$$\Pr[Z_1 = z, Z_{k+1} = z'] := \frac{\sum_{i=1}^{L-k} P_{Z_i Z_{i+k}}(z, z')}{L - k} =: (P_k^{ZZ'})_{zz'}. \quad (9)$$

Then, apply Eq. (1) again, we have the skip-grams for the generated and real sign cluster units satisfy

$$G^\top P_k^{XX'} G = P_k^{YY'}, \quad 1 \leq k \leq K - 1. \quad (10)$$

Again, we approximate this constraint by minimizing their ℓ_1 distance:

$$\mathcal{L}_{\text{skip}}(G) = \sum_{k=1}^K \|G^\top \hat{P}_k^{XX'} G - \hat{P}_k^{YY'}\|_1. \quad (11)$$

The overall loss for the word-level speech2sign-U is then

$$\mathcal{L}_{\text{sp2sign-U,word}} = \mathcal{L}_{\text{pos}} + \lambda \mathcal{L}_{\text{skip}}. \quad (12)$$

Speech-to-sign retriever Given a query speech audio (sign video), we would like to use it to retrieve its translation from a database of sign videos (speech audios). To this end, we use the generator

	# signs	# train sents	# valid	# test
Character-level datasets				
FS LibriSpeech	87k	287k	5.5k	5.6k
FS LJSpeech	87k	13.1k	348	523
Word-level datasets				
ASL Libri. 100	2.6k	14.1k	56	54
ASL Libri. 200	4.4k	56.2k	291	311
ASL Libri. 500	8.2k	137k	941	956
ASL Libri. 1k	11.6k	290k	2.7k	2.5k

Table 1: Dataset statistics

211 to compute a similarity score between each speech
212 sequence X and sign sequence Y as:

$$213 \quad \text{Sim}(X, Y) = -\frac{1}{L} \text{DTW}(G(X), Y), \quad (13)$$

214 where $\text{DTW}(\cdot, \cdot)$ is the dynamic time warping dis-
215 tance between two feature sequences with *cosine*
216 *distance* as the frame-level metric, computed using
217 the DTW library (Giorgino, 2009).

218 4 Experiments

219 4.1 Datasets

220 The detailed statistics are shown in Table 1.

221 **Fingerspelling LibriSpeech** To extract semantic
222 units from the fingerspelling signs, we trained
223 the visual CPC encoder on a sentence-level fin-
224 gerspelling dataset constructed from the 960-hour
225 LibriSpeech dataset and the Unvoiced dataset (Na-
226 garaj, 2018). To construct the dataset, we replace
227 each letter in the LibriSpeech transcript with an
228 image of that letter’s ASL Alphabet symbol chosen
229 uniformly at random from Unvoiced. To study the
230 effect of visual variability on SSR-U, we subset the
231 ASL Alphabet images to 100, 300, 500, or 1000
232 images per letter sign. The dev-clean subset of
233 LibriSpeech is used as the validation set.

234 **Fingerspelling LJSpeech** We train our character-
235 level model on another sentence-level finger-
236 spelling dataset constructed from LJSpeech (Ito
237 and Johnson, 2017) and the ASL Alphabet dataset
238 similar to the fingerspelling LibriSpeech.

239 **ASL LibriSpeech** For the word-level SSR-U,
240 we construct another corpus using LibriSpeech
241 for speech and MSASL (Joze and Koller, 2019)
242 for word-level sign videos. Since many MSASL
243 videos no longer exist on YouTube, only 11.6k out
244 of 25k videos are downloaded. Further, due to the

mismatch in vocabulary size, we use forced align-
245 ment information to filter out LibriSpeech words
246 that don’t appear in MSASL and keep sentences
247 that are at least 5 words long. Next, for each word
248 in each sentence, we pick a word-level sign video
249 uniformly at random from MSASL. To study the
250 effect of vocabulary size on our model, we follow
251 the split provided by (Joze and Koller, 2019) to
252 subset the data to a vocabulary size of 100, 200,
253 500 or 1000.

255 4.2 Overall results

Evaluation metrics We evaluate the perfor-
256 mance of our systems using two metrics: the *unit*
257 *error rate* (UER) is the average insertion I , dele-
258 tion D , and substitution S error between the pre-
259 dicted and true visual cluster units, which may be
260 character- or word-level units depending on the
261 task:

$$262 \quad \text{UER} = \frac{I + D + S}{3} \times 100.$$

The other metric we used to evaluate the speech-to-
264 sign ($A \rightarrow V$) and sign-to-speech ($V \rightarrow A$) retrieval
265 tasks is *recall@k* ($R@k$) ($k = 1, 5, 10$), which is
266 the percentage of hits in the top k results returned
267 by the retriever.

Character-level SSR-U The character-level re-
269 sults are shown in Table 2. To obtain retrieval
270 results, we trained our own wav2vec-U 2.0 using
271 the code released by the authors. Unfortunately,
272 we were unable to achieve the same results they
273 report in their paper. For our ASR-U experiments,
274 wav2vec-U significantly outperforms wav2vec-U
275 2.0 in terms of both word error rates and retrieval
276 tasks. For SSR-U, we compare our models with
277 wav2vec-U (and 2.0) as well as a supervised im-
278 age and caption retrieval model trained under a
279 ranking-based criterion (Harwath et al., 2018). We
280 replace their original CNN speech encoder with
281 a two-layer MLP with hidden and output sizes of
282 256 and ReLU activation, and their VGG16 image
283 encoder with a linear image encoder with an out-
284 put size of 256. We found that our models with
285 100 and 300 images per letter achieve superior per-
286 formances in terms of recall scores, even to the
287 text-based wav2vec-U, but remain about 30% be-
288 low the supervised topline. Notably, our model
289 performs worse on the $A \rightarrow V$ direction than on
290 the $V \rightarrow A$ direction, especially in terms of re-
291 call@1. This is perhaps due to significant insertion

Model	Images/ltr	UER↓	A→V			V→A		
			R@1↑	R@5↑	R@10↑	R@1↑	R@5↑	R@10↑
<i>Supervised speech-to-sign recognition</i>								
(Harwath et al., 2018)	1000	-	85.1	93.7	95.4	79.3	96.8	99.1
<i>Unsupervised speech recognition</i>								
wav2vec-U	-	39.5	1.9	42.1	59.5	17.6	44.2	65.2
wav2vec-U 2.0	-	68.1	1.1	6.3	11.8	2.6	9.8	14.9
<i>Unsupervised speech-to-sign recognition</i>								
speech2sign-U	100	43.1	1.7	42.1	63.4	27.5	62.7	78.4
speech2sign-U	300	45.0	1.9	48.8	67.1	22.8	53.9	71.5
speech2sign-U	500	46.2	1.0	33.1	57.2	31.3	57.0	72.8
speech2sign-U	1000	48.6	1.3	43.8	63.8	32.5	58.7	71.5

Table 2: Overall speech2sign-U results on FS LJSpeech

errors in the generated character sequence, which leads to many false positives during speech-to-sign retrieval.

Word-level SSR-U The word-level results are shown in Table 3. To establish top-line results for error rates and retrieval recall scores, we train a word-level unsupervised speech recognition model, speech2text-U, using the same criterion as speech2sign-U in Eq. 12, except by replacing the sign cluster sequences obtained from clustering word-level sign video features (see Section 3.2) with the underlying textual word labels as the target random variable Y . At the same time, for the subset with a vocabulary size of 98, we compare the performance of our model that uses unsupervised unigram and skipgram matching with wav2vec-U, which uses a JSD GAN for distribution matching, to show our proposed training method significantly improves the word error rates and the recall scores for both retrieval directions. However, we still observe a large gap in recall between our unsupervised model and the supervised speech-to-image retrieval model (Harwath et al., 2018). The performance of both word-level ASR-U and SSR-U degrades as the vocabulary size increases. The unit error rate (UER) increases from 53.6% to 87.9%, the recall@1 of speech-to-sign (A→V) retrieval decreases from 69.6% to 12.1%, and the recall@1 of sign-to-speech (V→A) retrieval decreases from 71.4% to 10.9% as the vocabulary size increases from 98 to 877. Such performance degradation is much more significant than that of character-level SSR-U because the word modality involves extra

morphological complexity on top of the phonological character modality.

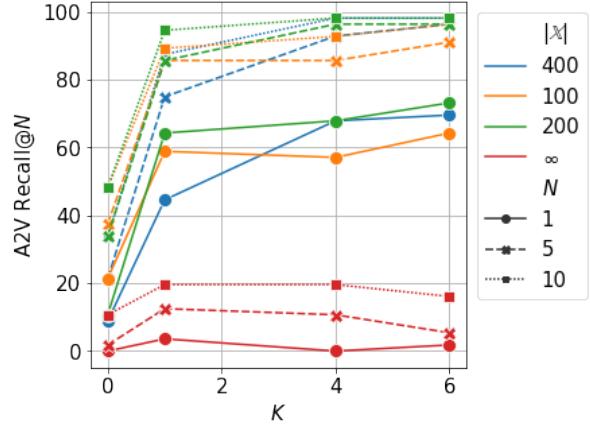


Figure 2: speech2sign-U retrieval results (recall@ N , $N \in \{1, 5, 10\}$) vs skip-gram size K for various number of speech clusters on ASL LibriSpeech 100

4.3 Analysis

Effect of skip-gram size The relation between recall@1, 5, 10 and skip-gram size K is shown in Figure 2. Increasing K generally improves all recall metrics for SSR-U by introducing more constraints to the generator mapping, though the performance starts to saturate at $K = 4$.

Effect of the number of speech clusters We experiment with speech2sign-U models with speech cluster sizes $|\mathbb{X}|$ equal to 100, 200, 400, and a model that directly takes raw wav2vec 2.0 features as inputs ($|\mathbb{X}| = \infty$), as shown in Figure 2. We found that the continuous model is significantly

Model	Vocab size	UER \downarrow	A \rightarrow V			V \rightarrow A		
			R@1 \uparrow	R@5 \uparrow	R@10 \uparrow	R@1 \uparrow	R@5 \uparrow	R@10 \uparrow
<i>Supervised speech-to-sign recognition</i>								
(Harwath et al., 2018)	877	-	55.2	78.9	86.3	51.7	85.5	93.1
<i>Unsupervised speech recognition</i>								
wav2vec-U	98	73.7	16.1	32.1	51.8	17.8	41.1	50.0
speech2text-U	98	7.5	98.2	98.2	100	98.2	98.2	98.2
speech2text-U	193	11.2	96.9	98.6	99.0	96.9	99.3	99.3
speech2text-U	468	30.0	68.0	86.2	90.2	66.7	85.3	90.2
speech2text-U	877	34.4	37.9	60.7	69.3	38.7	59.4	68.3
<i>Unsupervised speech-to-sign recognition</i>								
speech2sign-U	98	53.6	69.6	96.4	98.2	71.4	96.4	100
speech2sign-U	193	60.8	75.3	91.1	92.4	69.4	90.0	93.5
speech2sign-U	468	73.2	56.5	76.8	83.6	47.6	74.5	83.5
speech2sign-U	877	87.9	12.1	25.5	32.6	10.9	22.1	29.7

Table 3: Overall speech2sign-U results on ASL LibriSpeech

Video features	Vocab Size	A \rightarrow V		
		R@1	R@5	R@10
VGG19	98	32.1	64.3	78.6
OpenPose	98	0.0	8.9	16.1
	98	76.8	89.3	96.4
I3D RGB	193	66.0	86.3	91.4
	877	1.1	3.0	5.1
	98	69.6	96.4	98.2
I3D flow	193	63.9	86.9	91.1
	468	43.9	68.5	76.6
	877	12.1	25.5	32.6
	98	28.6	67.9	76.8
I3D joint	193	75.3	91.1	92.4
	468	56.5	76.8	83.6
	877	0.1	0.2	0.4

Table 4: Effect of the video features on ASL LibriSpeech with various vocabulary sizes

worse than discrete models and $|\mathbb{X}| = 200$ provides the most consistent recall scores across different skip-gram sizes.

Effect of training objectives The effect of different training objectives including the default speech2sign-U loss (L1) in Eq. (12), the maximum mean discrepancy (MMD) GAN and the Jensen-Shannon divergence (JSD) GAN is shown in Figure 3. For models trained with MMD and JSD

Boundary Label	Word Boundary F1	A \rightarrow V		
		R@1	R@5	R@10
speech2text-U				
Word	100	98.2	98.2	100
Phoneme	88.1	78.6	98.2	98.2
speech2sign-U				
Word	100	69.6	96.4	98.2
Phoneme	88.1	57.1	82.1	91.1

Table 5: Effect of the speech segmentation using speech2sign-U on ASL LibriSpeech 100

GAN loss, we instead feed the generator outputs to a discriminator with a single convolutional layer while keeping all other settings the same. Our experiment indicates that the GAN-free approach is consistently more stable and accurate compared to the GAN-based approach.

Effect of visual features The effect of visual features is shown in Table 4. We experimented with different types of visual features on ASL LibriSpeech with different vocabulary sizes such as VGG19 and the pose keypoint features from OpenPose (Cao et al., 2019). For the OpenPose features, we extract the keypoints from each video frames and re-sample each sign video feature frames to 30 frames as the segment-level feature. I3D architecture (Carreira and Zisserman, 2017) significantly

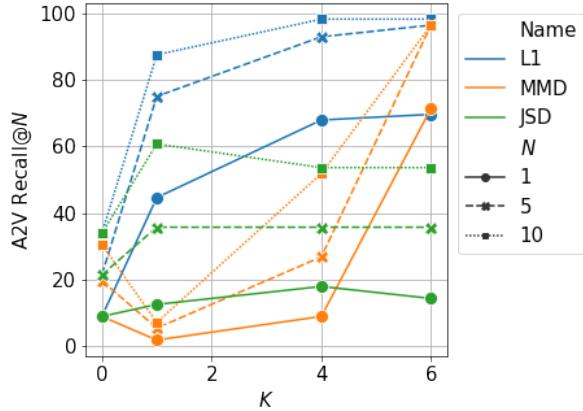


Figure 3: Retrieval results (recall@ N , $N \in \{1, 5, 10\}$) vs skip-gram size K for various types of training objective on ASL LibriSpeech 100

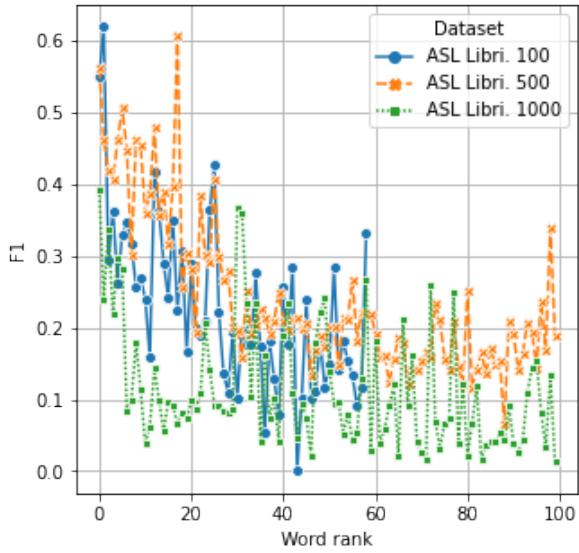


Figure 4: Word detection F1 vs. word rank by frequency

outperforms VGG19 and OpenPose as a feature extractor, demonstrating the importance of temporal information for SSR-U. We also found that I3D with optical flow features performs better than I3D with raw RGB inputs for most vocabulary sizes. Further, we found that concatenating the features from the RGB-based and flow-based I3Ds is beneficial for vocabulary sizes 193 and 468 but not when the vocabulary size is too small or too large, even causing training instability for vocabulary size 877.

Effect of segmentations The effect of gold and predicted speech segmentation for word-level SSR-U is shown in Table 5. For models trained with phoneme boundaries, we obtain predicted word segmentations using a CPC-based unsupervised segmentation system (Kreuk et al., 2020) with mean-

pooled phoneme-level wav2vec 2.0 features as inputs. The convolutional encoder in the original model is replaced by a two-layer MLP with 256 output dimensions trained on ASL LibriSpeech 100 for 200 epochs. This yields an exact-match boundary F1 of 88%. Using such detected word boundaries, we found about a 20% drop in recall@1 for speech2text-U and an 8-17% relative drop in recall@1,5,10 for speech2sign-U. Still, our model remains much better than the wav2vec-U baseline with ground-truth word boundaries, demonstrating its robustness to segmentation noise.

Effect of word frequencies We plotted the F1 score of the first 100 word classes ranked by frequency in Figure 4. For ASL LibriSpeech 100 and 500, while noisy, it is not hard to observe that the F1 score positively correlates with word frequency in a somewhat exponential fashion. Starting with F1 above 0.55 for the most frequent word, the performance quickly drops below 0.2 at around the 30th most frequent word. This trend is less conclusive on ASL LibriSpeech 1000 with generally low F1 scores, but the highest F1 scores are still observed for the most frequent words. The trend is also illustrated by the DTW alignment of a speech-video pair correctly retrieved by speech2sign-U in Figure 5. In our example, speech2sign-U mistakes the sign “more” for more frequent signs such as “when” and “have”. Additional factors such as *visual similarity* also play a role in the case of “more” and “when”, as both signs involve touching the tips of both hands. Such factors may explain the fluctuations in Figure 4. More error analysis can be found in Appendix A.

5 Related works

Sign language recognition One way to bridge between sign language and written/spoken language is to build a sign language recognition (SLR) system trained on parallel sign language and text corpora. The earliest attempts tried to recognize fingerspelling gestures using hand-tracking signals from wired gloves (Grimes, 1983; Charayaphan and Marble, 1992). Later works introduced vision to either correct the errors made by the hand-tracking model, or to serve as a cheaper and less-intrusive alternative (Tamura and Kawasaki, 1988). Focusing on the problem of *isolated* sign recognition and treating it as a classification task, a variety of statistical and deep learning models have been proposed, such as HMM (Starner and Pent-

366
367
368
369
370
371
372
373
374
375

376
377
378
379
380
381

382
383
384
385
386
387
388
389
390
391
392
393

394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415

416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

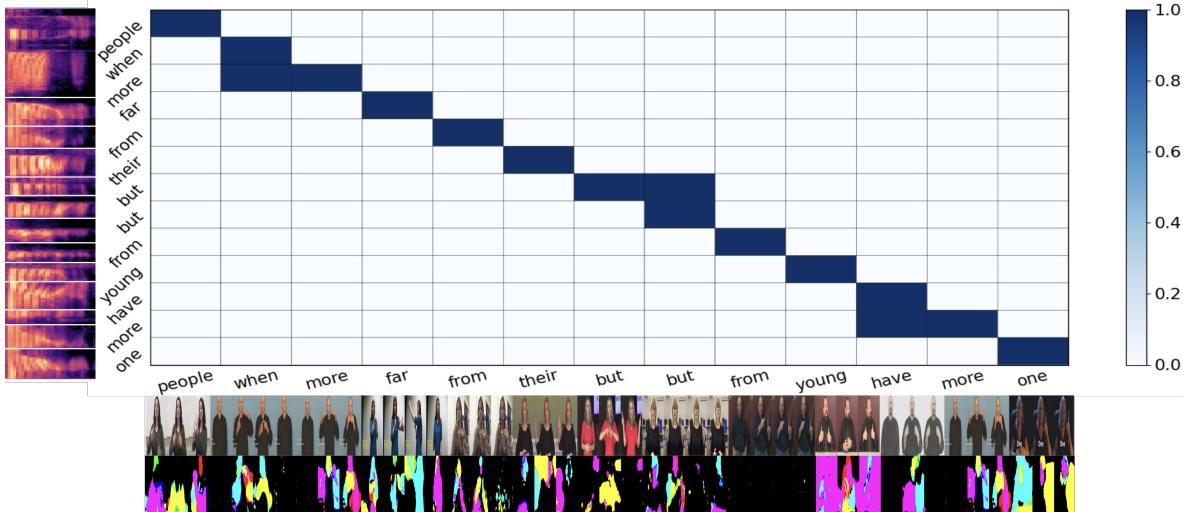


Figure 5: An example of the DTW alignment by speech2sign-U between a pair of speech and sign video (with its optical flow sequence shown below)

land, 1997), 3D-CNN (Huang et al., 2015), two-stream inflated 3D (I3D) CNN (Carreira and Zisserman, 2017; Joze and Koller, 2019), and transformer (Boháček and Hrúz, 2022), among others. To handle multi-sign video sequences, (Koller et al., 2016, 2017, 2018, 2019) reformulate the problem as a sequence labeling problem and develop various systems based on 2D-CNN-HMM hybrid models for German sign language recognition. Later works improve the alignment mechanism of previous models using soft DTW (Huang et al., 2018), CTC with DTW constraints (Pu et al., 2019) or pseudo-labeling refinement (Zhou et al., 2019). While some aim to directly use raw RGB images or generic action features like optical flow as inputs (Koller et al., 2016; Huang et al., 2018; Joze and Koller, 2019), others have found domain-specific features like whole-body and hand keypoints to be more reliable and robust (Boháček and Hrúz, 2022). Thanks to the rapid development of the field, there are now many word-level and sentence-level datasets available in different SLs, and we refer to (Joze and Koller, 2019) for a more comprehensive review.

Unsupervised cross-modal alignment The task of translating between two languages without parallel corpora has been demonstrated between written language pairs (MT-U) and between spoken-written language pairs (ASR-U). (Haghghi et al., 2008) and (Ravi and Knight, 2011; Pourdamghani and Knight, 2017) are respectively the first to treat word-level and sentence-level MT-U as a distribution matching problem and built the first

such systems by training statistical machine translation systems using nonparallel corpora, which are further improved by (Artetxe et al., 2018b). To allow more general source and target distributions, (Zhang et al., 2017a,b; Conneau et al., 2018; Artetxe et al., 2018a; Lample et al., 2018) instead use neural networks to embed the source and target distributions and match the distributions using either shared denoising autoencoder (Artetxe et al., 2018a), earth-mover distance minimization (Zhang et al., 2017b) or a generative adversarial network (GAN) with additional regularization losses (Zhang et al., 2017a; Conneau et al., 2018; Lample et al., 2018). (Chung et al., 2018; Liu et al., 2018; Chen et al., 2019; Baevski et al., 2021; Liu et al., 2022) adapt and perfect the GAN-based approach for spoken-written language pairs by leveraging large-scale self-supervised speech representation learning models (Chung and Glass, 2018; Baevski et al., 2020) as well as iterative self-training techniques (Liu et al., 2018).

6 Conclusion

In this paper, we propose the task of unsupervised speech-to-sign language recognition and a neural network model, speech2sign-U, capable of both character-level and word-level SSR-U. On various unpaired speech and ASL datasets, our models consistently outperform previous unsupervised models such as wav2vec-U. Further, we found our model reliable to train for a variety of vocabulary sizes and robust against various types of noise in both speech and visual modalities.

497 7 Limitations

498 Our model currently requires high-quality word
499 boundaries for both speech and sign videos. How-
500 ever, as demonstrated by our preliminary results in
501 Table 5, we can overcome such limitations by incor-
502 porating more powerful unsupervised segmentation
503 algorithms to our system. Further, while our dataset
504 is sufficient to model the variability in speech and
505 videos, all experiments to date have assumed that
506 spoken and signed sentences share similar word
507 order, which may not be true of natural spoken
508 and signed communications. A future direction
509 of this research will seek to develop methods for
510 spoken-sign language pairs with very different syn-
511 tactic structures. Lastly, the vocabulary size under
512 our study on word-level SSR-U is relatively small
513 (<1000), and a promising future direction is to ex-
514 tend the current approach to deal with much larger
515 vocabulary size in more diverse conversations.

516 8 Ethical considerations

517 One potential ethical concern for our model is
518 the risk of miscommunication. Due to the small
519 amount of resources used to train our system, it
520 tends to be less accurate than its supervised coun-
521 terpart, and its mistakes may cause confusion, mis-
522 understanding and other psychological harm to the
523 users of our systems. The other ethical concern is
524 that the data used to train the system is demographi-
525 cally homogeneous, as we have noticed from some
526 brief inspections that most of the signers in the
527 ASL datasets are white middle-aged adults. This
528 may lead the system to worse retrieval accuracy
529 for people underrepresented in the training corpus,
530 such as black people, children and elderly people.

531 References

- 532 Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018a.
533 [Unsupervised neural machine translation](#). In *International Conference on Learning Representations*.
- 535 Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018b.
536 [Unsupervised statistical machine translation](#). In *Proceedings of the 2018 Conference on Empirical Meth-
537 ods in Natural Language Processing*, pages 3632–
538 3642, Brussels, Belgium. Association for Computa-
539 tional Linguistics.
- 541 Alexei Baevski, Wei-Ning Hsu, Alexis Conneau, and
542 Michael Auli. 2021. [Unsupervised speech recogni-
543 tion](#). In *Neural Information Processing System*.
- 544 Alexei Baevski, Henry Zhou, Abdelrahman Mohamed,
545 and Michael Auli. 2020. [wav2vec 2.0: A framework](#)

- 546 for self-supervised learning of speech representations.
547 In *Neural Information Processing System*.
- 548 S. Bhati, J. Villalba, P. Želasko, L. Moro-Velázquez,
549 and N. Dehak. 2021. [Segmental contrastive predic-
550 tive coding for unsupervised word segmentation](#). In
551 *Interspeech*, page 366–370.
- 552 Matyáš Boháček and Marek Hrúz. 2022. [Sign pose-
553 based transformer for word-level sign language recog-
554 nition](#). In *Proceedings of the IEEE/CVF Winter Con-
555 ference on Applications of Computer Vision (WACV)
556 Workshops*.
- 557 Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei,
558 and Yaser Sheikh. 2019. [Openpose: Realtime multi-
559 person 2d pose estimation using part affinity fields](#).
560 *IEEE Transactions on Pattern Analysis and Machine
561 Intelligence*.
- 562 Joao Carreira and Andrew Zisserman. 2017. [Quo vadis,
563 action recognition? a new model and the kinetics
564 dataset](#). In *The IEEE / CVF Computer Vision and
565 Pattern Recognition Conference (CVPR)*.
- 566 C. Charayaphan and A. E. Marble. 1992. [Image pro-
567 cessing system for interpreting motion in american
568 sign language](#). *Journal of Biomedical Engineering*,
569 14(5):419–425.
- 570 Kuan-Yu Chen, Che-Ping Tsai, Da-Rong Liu, Hung-Yi
571 Lee, and Lin shan Lee. 2019. [Completely unsuper-
572 vised speech recognition by a generative adversarial
573 network harmonized with iteratively refined hidden
574 markov models](#). In *Interspeech*.
- 575 Yu-An Chung and James Glass. 2018. [Speech2vec: A
576 sequence-to-sequence framework for learning word
577 embeddings from speech](#). In *INTERSPEECH*.
- 578 Yu-An Chung, Wei-Hung Weng, Schrasling Tong, and
579 James Glass. 2018. [Unsupervised cross-modal align-
580 ment of speech and text embedding spaces](#). In *Neural
581 Information Processing System*.
- 582 A. Conneau, G. Lample, M. Ranzato, L. Denoyer, and
583 H. Jégou. 2018. [Word translation without parallel
584 data](#). In *International Conference on Learning Rep-
585 resentations*.
- 586 Santiago Cuervo, Adrian Lańcucki, Ricard Marxer,
587 Paweł Rychlikowski, and Jan Chorowski. 2022.
588 [Variable-rate hierarchical cpc leads to acoustic unit
589 discovery in speech](#). In *Neural Information Process-
590 ing System*.
- 591 Toni Giorgino. 2009. [Computing and visualizing dy-
592 namic time warping alignments in r: The dtw pack-
593 age](#). *Journal of Statistical Software*, 31(7):1–24.
- 594 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza,
595 Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron
596 Courville, and Yoshua Bengio. 2014. [Generative
597 adversarial nets](#). In *Neural Information Processing
598 System*.

599	Gary J. Grimes. 1983. <i>Digital data entry glove interface device</i> . US Patent.	600	
601	Aria Haghghi, Percy Liang, Taylor Berg-Kirkpatrick, and Dan Klein. 2008. <i>Learning bilingual lexicons from monolingual corpora</i> . In <i>Proceedings of ACL-08: HLT</i> , pages 771–779, Columbus, Ohio. Association for Computational Linguistics.	602	
603		604	
605		606	
607	David Harwath, Adrià Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass. 2018. <i>Jointly discovering visual objects and spoken words from raw sensory input</i> . In <i>ECCV</i> .	608	
609		610	
611	Jie Huang, Houqiang Li Wengang Zhou, and Weiping Li. 2015. Sign language recognition using 3d convolutional neural networks. In <i>IEEE International Conference on Multimedia and Expo (ICME)</i> , pages 612 1–6.	613	
614		615	
616	Jie Huang, Wengang Zhou, Qilin Zhang, Houqiang Li, and Weiping Li. 2018. <i>Video-based sign language 617 recognition without temporal segmentation</i> . In <i>AAAI 618 Conference on Artificial Intelligence (AAAI)</i> .	619	
620		621	
622	Keith Ito and Linda Johnson. 2017. <i>The lj 623 speech dataset</i> . https://keithito.com/LJ-Speech-Dataset/ .	624	
625		626	
627	Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. 628 Billion-scale similarity search with GPUs. <i>IEEE 629 Transactions on Big Data</i> , 7(3):535–547.	630	
631		632	
633	Hamid Vaezi Joze and Oscar Koller. 2019. <i>MS-ASL: A 634 large-scale data set and benchmark for understanding 635 american sign language</i> . In <i>The British Machine 636 Vision Conference (BMVC)</i> .	637	
638		639	
640	Diederik P. Kingma and Jimmy Lei Ba. 2015. <i>Adam: A 641 method for stochastic optimization</i> . In <i>International 642 Conference on Learning Representations</i> .	643	
644		645	
646	Oscar Koller, Necati Cihan Camgoz, Hermann Ney, and 647 Richard Bowden. 2019. Weakly supervised learning with multi-stream CNN-LSTM-HMMs to discover sequential parallelism in sign language videos. <i>IEEE 648 Transactions on Pattern Analysis and Machine Intelligence</i> .	649	
650		651	
652	Oscar Koller, Sepehr Zargaran, and Hermann Ney. 2017. <i>Re-sign: Re-aligned end-to-end sequence modelling with deep recurrent CNN-HMMs</i> . In <i>The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)</i> , pages 4297–4305.	653	
654		655	
656	Oscar Koller, Sepehr Zargaran, Hermann Ney, and 657 Richard Bowden. 2016. <i>Deep sign: Hybrid CNN-HMM for continuous sign language recognition</i> . In <i>Proc. British Machine Vision Conference (BMVC)</i> , page 1–12.	658	
659		660	
661	Oscar Koller, Sepehr Zargaran, Hermann Ney, and 662 Richard Bowden. 2018. <i>Deep sign: Enabling robust statistical continuous sign language recognition via hybrid CNN-HMMs</i> . <i>International Journal of 663 Computer Vision (IJCV)</i> , 126(12):1311–1325.	664	
665		666	
667		668	
669	Felix Kreuk, Joseph Keshet, and Yossi Adi. 2020. <i>Self-supervised contrastive learning for unsupervised phoneme segmentation</i> . In <i>INTERSPEECH</i> .	670	
671		672	
673	Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018. <i>Unsupervised machine translation using monolingual corpora only</i> . In <i>International Conference on Learning Representations</i> .	674	
675		676	
676	Alexander H. Liu, Wei-Ning Hsu, Michael Auli, and Alexei Baevski. 2022. <i>Towards end-to-end unsupervised speech recognition</i> . In <i>ArKiv</i> .	677	
678		679	
680	Da-Rong Liu, Kuan-Yu Chen, Hung-Yi Lee, and Lin 681 shan Lee. 2018. Completely unsupervised phoneme 682 recognition by adversarially learning mapping relationships 683 from audio embeddings. In <i>Interspeech</i> .	684	
685		686	
686	Akash Nagaraj. 2018. <i>ASL alphabet</i> .	687	
688		689	
690	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. <i>Representation learning with contrastive predictive 691 coding</i> . In <i>ArKiv</i> .	692	
693		694	
695	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael 696 Auli. 2019. <i>fairseq: A fast, extensible toolkit for 697 sequence modeling</i> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association 698 for Computational Linguistics (Demonstrations)</i> , 699 pages 48–53, Minneapolis, Minnesota. Association 700 for Computational Linguistics.	701	
702		703	
704	Adam Paszke, Sam Gross, Francisco Massa, Adam 705 Lerer, James Bradbury, Trevor Gregory, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. <i>Pytorch: An imperative style, high-performance deep learning library</i> . In <i>Advances in Neural Information Processing Systems 32</i> , pages 8024–8035. Curran Associates, Inc.	706	
707		708	
709	Nima Pourdamghani and Kevin Knight. 2017. <i>Deciphering 710 related languages</i> . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2513–2518, Copenhagen, Denmark. Association for Computational Linguistics.	711	
712		713	
714	Junfu Pu, Wengang Zhou, and Houqiang Li. 2019. <i>Iterative alignment network for continuous sign language 715 recognition</i> . In <i>The IEEE / CVF Computer Vision and 716 Pattern Recognition Conference (CVPR)</i> .	717	
718		719	
720	Sujith Ravi and Kevin Knight. 2011. <i>Deciphering foreign 721 language</i> . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 12–21, Portland, Oregon, USA. Association for Computational Linguistics.	722	
723		724	
725		726	
727		728	
729		730	
731		732	
733		734	
735		736	
737		738	
739		740	
741		742	
743		744	
745		746	
747		748	
749		750	
751		752	
753		754	
755		756	
757		758	
759		760	
761		762	
763		764	
765		766	
767		768	
769		770	
771		772	
773		774	
775		776	
777		778	
779		780	
781		782	
783		784	
785		786	
787		788	
789		790	
791		792	
793		794	
795		796	
797		798	
799		800	
801		802	
803		804	
805		806	
807		808	
809		810	
811		812	
813		814	
815		816	
817		818	
819		820	
821		822	
823		824	
825		826	
827		828	
829		830	
831		832	
833		834	
835		836	
837		838	
839		840	
841		842	
843		844	
845		846	
847		848	
849		850	
851		852	
853		854	
855		856	
857		858	
859		860	
861		862	
863		864	
865		866	
867		868	
869		870	
871		872	
873		874	
875		876	
877		878	
879		880	
881		882	
883		884	
885		886	
887		888	
889		890	
891		892	
893		894	
895		896	
897		898	
899		900	
901		902	
903		904	
905		906	
907		908	
909		910	
911		912	
913		914	
915		916	
917		918	
919		920	
921		922	
923		924	
925		926	
927		928	
929		930	
931		932	
933		934	
935		936	
937		938	
939		940	
941		942	
943		944	
945		946	
947		948	
949		950	
951		952	
953		954	
955		956	
957		958	
959		960	
961		962	
963		964	
965		966	
967		968	
969		970	
971		972	
973		974	
975		976	
977		978	
979		980	
981		982	
983		984	
985		986	
987		988	
989		990	
991		992	
993		994	
995		996	
997		998	
999		999	

706 Shaoqing Ren, Kaiming He, Ross Girshick, and Jian
707 Sun. 2015. Faster R-CNN: Towards real-time object
708 detection with region proposal networks. In *Neural*
709 *Information Processing System*.

710 Thad Starner and Alex Pentland. 1997. Real-time Amer-
711 ican sign language recognition from video using
712 hidden markov models. *Motion-Based Recognition*,
713 page 227–243.

714 Shinichi Tamura and Shingo Kawasaki. 1988. *Recog-
715 nition of sign language motion images*. *Pattern Recog-
716 nition*, 21(4):343–353.

717 Meng Zhang, Yang Liu, Huanbo Luan, and Maosong
718 Sun. 2017a. *Adversarial training for unsupervised
719 bilingual lexicon induction*. In *Proceedings of the
720 55th Annual Meeting of the Association for Com-
721 putational Linguistics (Volume 1: Long Papers)*,
722 pages 1959–1970, Vancouver, Canada. Association
723 for Computational Linguistics.

724 Meng Zhang, Yang Liu, Huanbo Luan, and Maosong
725 Sun. 2017b. *Earth mover’s distance minimization for
726 unsupervised bilingual lexicon induction*. In *Proceed-
727 ings of the 2017 Conference on Empirical Methods
728 in Natural Language Processing*, pages 1934–1945,
729 Copenhagen, Denmark. Association for Compu-
730 tational Linguistics.

731 Hao Zhou, Wengang Zhou, and Houqiang Li. 2019.
732 Dynamic pseudo label decoding for continuous sign
733 language recognition. In *International Conference
734 on Multimedia and Expo (ICME)*.

735 A Appendix

736 A.1 Reproducibility checklist

737 All experiments are done on four 16GB NVIDIA
738 V100 GPUs and all models are implemented using
739 Pytorch (Paszke et al., 2019) and Fairseq (Ott et al.,
740 2019).

741 **Character-level speech2sign-U** We use the ex-
742 act same generator and discriminator architectures
743 as the wav2vec-U (Baevski et al., 2021). For the
744 CPC-based fingerspelling feature extractor, we use
745 a two-layer MLP as the encoder, with 256 hidden
746 units, ReLU activation and 256 output units and
747 a single-layer LSTM with 256 hidden and output
748 units as the autoregressive predictor. We found 3
749 prediction steps and 32 negative samples per pos-
750 itive sample for the CPC loss to be the best setting
751 for training. For the CPC-based fingerspelling fea-
752 ture extractor, we train for 60 epochs using Adam
753 optimizer (Kingma and Ba, 2015) with an initial
754 learning rate of 0.001, a batch size of 16 with
755 $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The checkpoint with

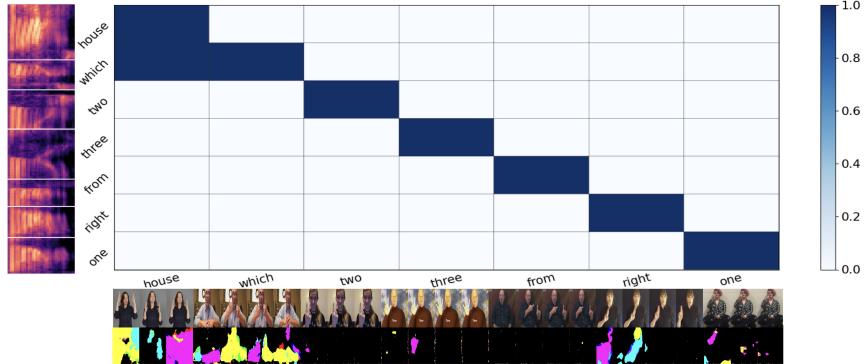
756 the highest average next-frame prediction perfor-
757 mance during training is used for the feature ex-
758 traction later. For the K-means clustering, we use
759 FAISS (Johnson et al., 2019) and set the number of
760 clusters to be the same as the vocabulary size. For
761 the GAN training, we train the model for 10000
762 updates and validate the model every 1000 updates
763 using the UER metric. We observe similar perfor-
764 mance between the best and the last checkpoints for
765 most experiments. Again, we follow the publicly
766 available implementation of wav2vec-U (Baevski
767 et al., 2021) using Fairseq for all the distributed
768 training, optimizer and scheduler setting.

769 **Word-level speech2sign-U** For extracting the
770 optical flow features of sign images, we use the
771 OpenCV implementation of Dual TV-L1 method
772 and resized all images to 224×224 . For the Open-
773 Pose features, we follow the default settings to
774 extract the pose keypoints and set the keypoint co-
775 ordinates to 0 when the model fails to detect any
776 keypoints. We also normalize the keypoints by the
777 size of the video frame. The I3D model we use
778 are trained on the ImageNet dataset and fine-tuned
779 on the Charades dataset, for both RGB and flow
780 implementations. The same CPC sign encoder as
781 that in character-level experiments is used, except
782 with the pretrained video features as inputs and the
783 outputs of the *MLP encoder* as outputs instead of
784 that of the LSTM model. We then train the CPC
785 sign encoder for 200 epochs on ASL LibriSpeech
786 1000. The CPC sign encoder features are then
787 quantized into the same number of discrete units as
788 the vocabulary size (100 for ASL LibriSpeech 100,
789 etc.) using K-means implemented in FAISS (John-
790 son et al., 2019). For the speech feature clustering,
791 we again use the FAISS (Johnson et al., 2019) im-
792 plementation of K-means with a cluster size of
793 about 4 times of the vocabulary size of ASL Libri-
794 Speech 100, 200 and 500, and 2000 clusters for
795 ASL LibriSpeech 1000. The cluster sizes are cho-
796 sen to ensure a cluster purity of about 90%. For
797 the word-level speech2sign-U, the speech genera-
798 tor is a linear layer with no bias. Skip grams of
799 a maximal step of 6 are used for experiments on
800 ASL LibriSpeech 100, 200 and 500, and a maximal
801 step of 4 are used for ASL LibriSpeech 1000. For
802 the unsupervised training, we train the model for
803 a number of updates equal to $3000 \times \left\lfloor \frac{\text{sample size}}{\text{batch size}} \right\rfloor$.
804 We found that larger batch size generally leads to
805 better performance, and use a batch size of 16k for
806 ASL LibriSpeech 100, 200 and 500, and a batch

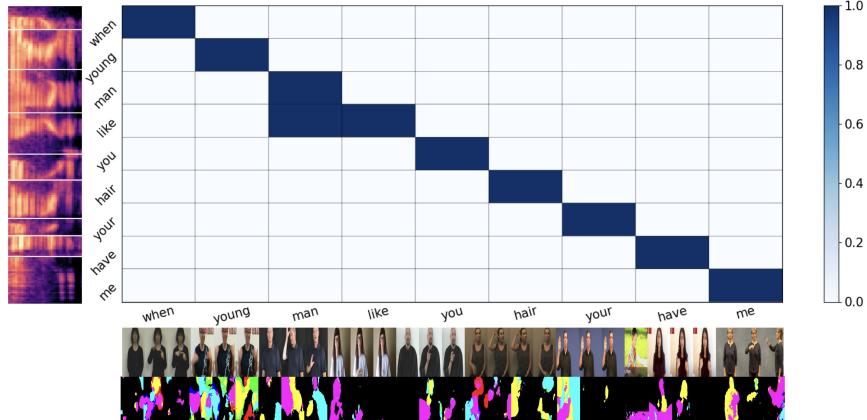
size of 12k for ASL LibriSpeech 1000 due to GPU memory constraints. Adam optimizer with a initial learning rate of 0.4 and $[\beta_1, \beta_2] = [0.9, 0.999]$ is used throughout the training.

A.2 More SSR-U retrieval examples and error analysis

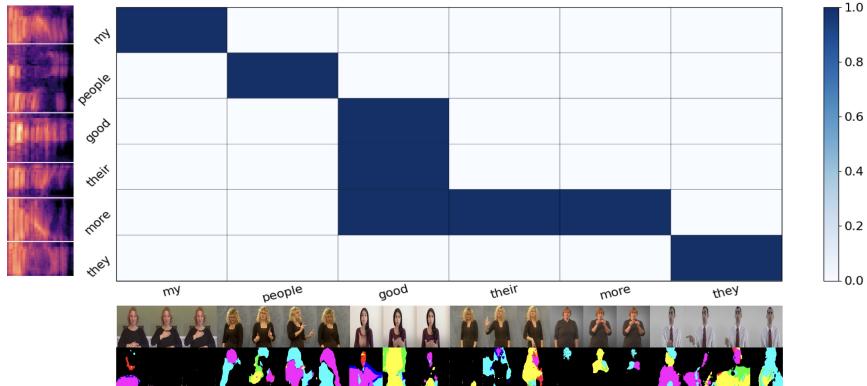
More DTW alignments between speech-video pairs correctly retrieved by speech2sign-U are shown in Figure 6. As we can see, our model is able to correctly align the speech and sign video after the DTW step. However, in order to better understand the type of errors the model is susceptible to, we also show the similarity map *before* the DTW step in Figure 7. While the similarity maps are noisier than their corresponding DTW alignments, the high similarity regions are correctly concentrated approximately along the diagonal most of the time. There are, however, several common failure modes by speech2sign-U. The most common mistake by the model is to confuse less frequent words with more frequent ones, for example, confuse the less frequent word “history” with the more frequent word “from” and “outside” in Figure 7d, or the less frequent “more” with the more frequent “good” in Figure 7c or the less frequent “like” with the more frequent “when” and “man” in Figure 7b. Another type of mistake is to confuse visually similar signs such as “one”, “two” and “three” in Figure 7a. The last common type of mistake for speech2sign-U is to confuse acoustically similar words, such as the word “they” and “their” in Figure 7c.



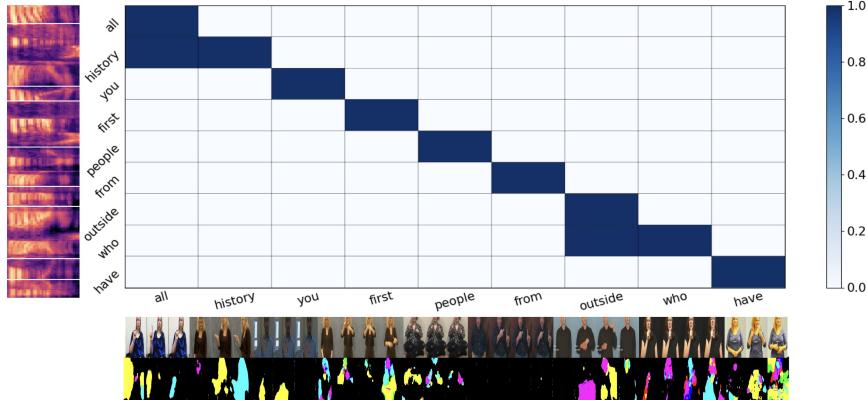
(a)



(b)



(c)



(d)

Figure 6: DTW alignments from ASL LibriSpeech 500

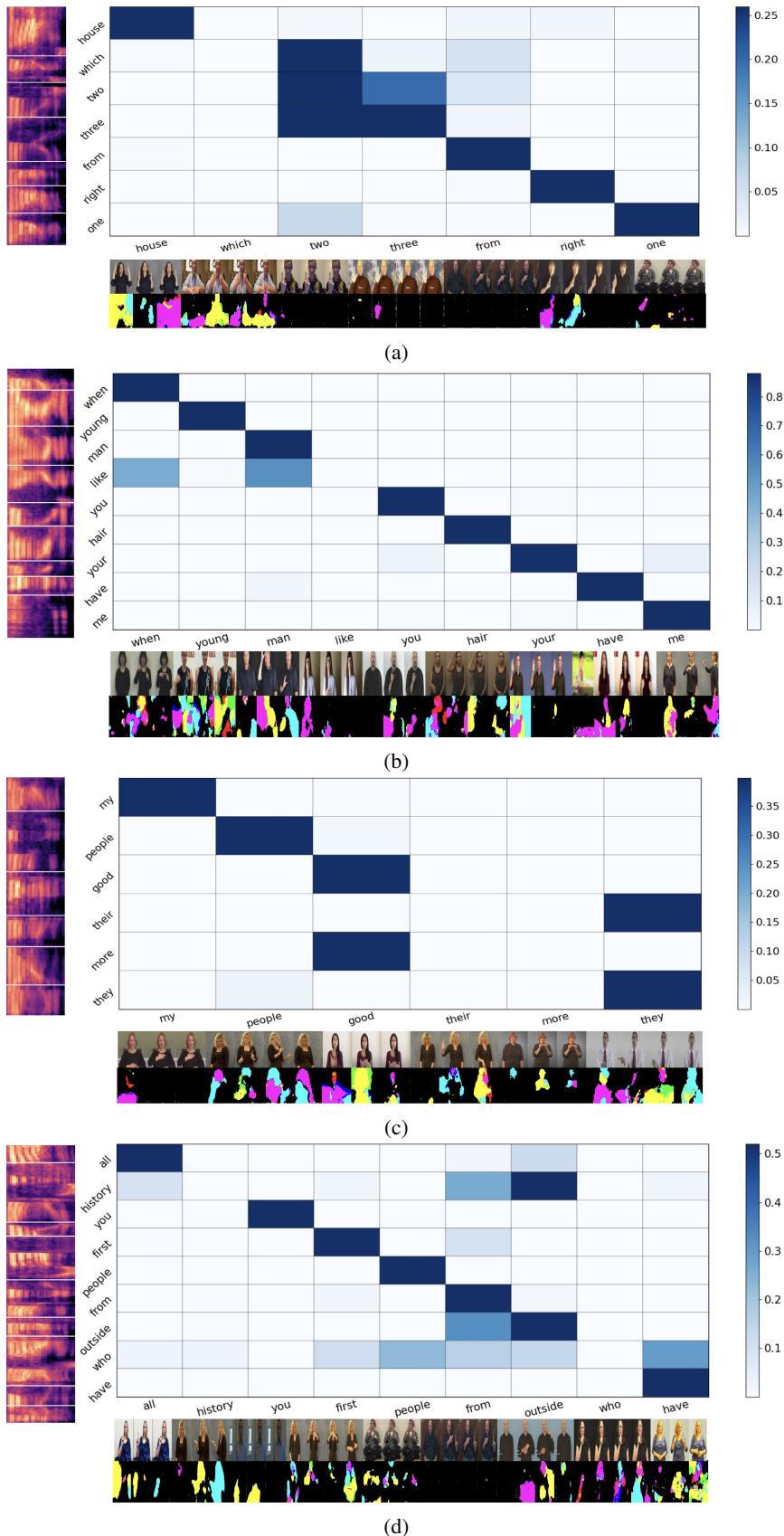


Figure 7: Similarity maps from ASL LibriSpeech 500