


# SiLVERScore: Semantically-Aware Embeddings for Sign Language Generation Evaluation

Saki Imai      Mert İnan      Anthony Sicilia      Malihe Alikhani

Northeastern University, Boston MA

{imai.s, m.alikhani}@northeastern.edu

## Abstract

Evaluating sign language generation is often done through back-translation, where generated signs are first recognized back to text and then compared to a reference using text-based metrics. However, this two-step evaluation pipeline introduces ambiguity: it not only fails to capture the multimodal nature of sign language—such as facial expressions, spatial grammar, and prosody—but also makes it hard to pinpoint whether evaluation errors come from sign generation model or the translation system used to assess it. In this work, we propose  SiLVERSCORE, a novel semantically-aware embedding-based evaluation metric that assesses sign language generation in a joint embedding space. Our contributions include: (1) identifying limitations of existing metrics, (2) introducing SiLVERScore for semantically-aware evaluation, (3) demonstrating its robustness to semantic and prosodic variations, and (4) exploring generalization challenges across datasets. On PHOENIX-14T and CSL-Daily datasets, SiLVERScore achieves near-perfect discrimination between correct and random pairs (ROC AUC = 0.99, overlap < 7%), substantially outperforming traditional metrics<sup>1</sup>.

## 1 Introduction

The ability to automatically evaluate sign language generation is critical for advancing accessibility and inclusion for the Deaf and Hard-of-Hearing (DHH) community, where collecting large scale human judgments remains expensive and challenging (Bragg et al., 2019; Yin et al., 2021; Huenerfauth et al., 2008). Scalable and reliable evaluation is necessary to ensure that generated sign language content meets user needs. Yet, progress toward fully automated systems is hindered by the absence of effective evaluation methods (Liu et al., 2023).

<sup>1</sup><https://github.com/sakimai/silverscore>

## Traditional Back-translation Metrics versus SiLVERScore

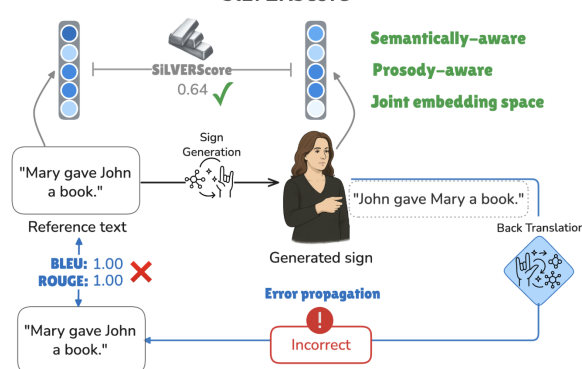




Figure 1: In this example, a sign language generation model accidentally swaps the referents, generating the sign for “John gave Mary a book” instead of “Mary gave John a book.” Traditional text-based metrics like BLEU and ROUGE (bottom left) rely on back-translation and fail to catch the error, assigning a perfect score because the English output matches the reference text, even though the meaning is incorrect. In contrast,  SiLVERSCORE (top) compares the generated signing video directly with the reference using a joint embedding space, correctly identifying the error.

To ensure outputs align with human expectations, we need robust evaluation metrics explicitly designed for the linguistic nature of sign language.

Automatically evaluating generated sign language remains challenging due to its unique multimodal linguistic nature, which incorporates facial expressions, manual markers, and spatiotemporal relationships into its prosody, iconicity, semantics, and pragmatics (Sandler, 2012; Liddell, 2003; Huenerfauth et al., 2008). Current evaluation methods rely on back-translation from visual to textual representations, which misaligns with the visual nature of sign language and leads to inaccuracies. As illustrated in Figure 1, a back-translation metric can assign a perfect BLEU or ROUGE score

even when the directional verb motion in the signing reverses the intended meaning. This mismatch motivates an evaluation paradigm that inspects the video itself rather than its textual translation.

While embedding-based metrics such as BLEURT (Sellam et al., 2020), BERTScore (Zhang\* et al., 2020) and CLIPScore (Hessel et al., 2021), have shown success in natural language processing, they have been underexplored for sign language evaluation. This limitation is primarily due to the scarcity and domain specificity of sign language datasets, which restrict the generalizability of sign embeddings. We hypothesize that these data limitations have hindered the development of effective embedding-based metrics for sign language generation.

To address this gap, we introduce  **SiLVERScore** (**Sign Language Video Embedding Representation Score**), a novel embedding-based metric for evaluating sign language generation. SiLVERScore directly compares generated and reference signs within a joint embedding space, capturing semantic and prosodic features.

Rather than asking whether embedding-based metrics are simply better than back-translation, our work investigates: *how embedding-based evaluation can more faithfully capture the linguistic and prosodic nuances of sign language*, and under what conditions it offers robust and generalizable performance. Our contributions are as follows:

1. We survey existing evaluation metrics for sign language generation and highlight their limitations (§ 2).
2. We introduce SiLVERScore, a novel semantically-aware embedding-based metric for evaluating sign language generation in a joint embedding space (§ 3).
3. We conduct prosodic and semantic tests to demonstrate that SiLVERScore outperforms traditional metrics (§ 4.2, § 4.3).
4. We perform a case study on generalization, the challenges of applying sign language models across different datasets and domains (§ 5).

## 2 Survey of Evaluation Metrics for Sign Language Processing

The evaluation of sign language generation systems has traditionally relied on back-translation approaches, first introduced by Camgoz et al. (2018). In these methods, a sign language translation model

(typically trained by the authors) is used to convert the generated signs into text for evaluation. However, the absence of a standardized sign-to-text translation system complicates this approach, introducing unknown error propagation.

To address these issues, researchers have proposed several multimodal metrics. For instance, Dynamic Time Warping Mean Joint Error (Huang et al., 2021) aligns generated and ground truth poses to measure spatial-temporal accuracy and compute the mean joint error. While effective for motion similarity, it penalizes valid linguistic variations that differ in pose but maintain semantic meaning. Similarly, Fréchet Gesture Distance (Yoon et al., 2020), Fréchet Video Distance (Unterthiner et al., 2019), Fréchet Inception Distance (Heusel et al., 2017) compare gesture distributions but focus on physical similarity rather than semantics (Hwang et al., 2022; Xie et al., 2024; Hwang et al., 2024; Dong et al., 2024). Common video quality scores (SSIM, PSNR, Inception Score, Temporal Consistency Metric) measure image quality, diversity, or smoothness, ignoring whether the sign is linguistically correct (Natarajan et al., 2022).

In a visual-spatial SignWriting domain, signwriting-evaluation (Moryossef et al., 2024) was proposed as a metric designed for this by using its symbol distance metric using the Hungarian algorithm (Kuhn, 1955). A sign language translation metric, SignBLEU (Kim et al., 2024) aims to mitigate the significant information loss due to the simplification to a single sequence of text for evaluation. However, despite its improvements, both remain confined to the text-realm.

Multimodal embedding-based methods are promising due to their ability to capture multimodal elements and eliminate errors introduced by back-translation. Existing sign language embeddings, such as SignCLIP (Jiang et al., 2024), offer a foundation for embedding-based evaluation. However, their application to sign language generation evaluation remains limited, primarily due to challenges in generalizability (§ 5). This paper aims to bridge this gap by introducing and validating a semantically aware embedding-based evaluation metric tailored to sign language generation.

## 3 SiLVERScore

The objective of SiLVERScore is to evaluate generated sign language videos without requiring a reference video. This evaluation measures the align-

ment between a sign video and its corresponding text by comparing their similarity in a shared joint embedding space, trained to capture multimodal relationships. The similarities are computed using CiCo (Cheng et al., 2023), a model that leverages contrastive learning to align video and text representations. This approach addresses the alignment issues discussed in § 5 by using a sliding window mechanism to localize alignment between modalities.

We employ CiCo due to its framework that: (i) formulates sign language retrieval as a cross-lingual retrieval task; (ii) demonstrates state-of-the-art performance on benchmarks such as PHOENIX-14T, CSL-Daily, and How2Sign; (iii) avoids reliance on pose estimation tools, eliminating dependency on pose extraction quality; and (iv) provides accessible code for implementation.

**Model Details** The sign encoder processes sign videos using a sliding window mechanism to generate embeddings. This approach enables the model to handle continuous video streams without requiring explicit segmentation at test time. This encoder combines domain-agnostic features, captured by a pre-trained I3D network (Varol et al., 2021) on BSL-1K, with domain-aware features from the same network fine-tuned on PHOENIX-14T/CSL-Daily. The features are weighted and fused before being fed into a 12-layer Transformer initialized with CLIP’s ViT-B encoder. The corresponding text is lowercased, byte pair encoded, and translated into English using Google Translate to align with the CLIP pretraining. The video and text embeddings are aligned through a contrastive learning objective with the InfoNCE loss. CiCo aligns video and text embeddings through a contrastive learning objective based on InfoNCE loss, which maximizes the similarity of matched video-text pairs while minimizing the similarity of unmatched pairs. This alignment is performed both globally across entire videos and texts and locally by retaining fine-grained mappings between video segments and individual text tokens. The resulting embeddings represent a semantically and temporally aware shared space that effectively captures the relationships between sign videos and their corresponding text annotations.

**Global Similarity Calculation** Global similarity is derived from a fine-grained similarity matrix

$$E \in R^{M \times L};$$

$$E(i, j) = S_i \cdot W_j^T, \quad (1)$$

where  $S_i \in R^D$  and  $W_j \in R^D$  represent video clip and word embeddings, respectively. To emphasize similarities, softmax re-weighting is applied:

$$E'(i, j) = \text{Softmax}(E(i, j)) \cdot E(i, j). \quad (2)$$

Row-wise summation followed by averaging yields the video-to-text similarity  $Z_{V2T}$ , while column-wise operations yield the text-to-video similarity  $Z_{T2V}$ .

In the implementation, the  $Z_{V2T}$  and  $Z_{T2V}$  similarities are equally weighted in the loss function. This equal weighting ensures that the global alignment of video-to-text and text-to-video pairs is equivalent, making it sufficient to use either  $Z_{V2T}$  or  $Z_{T2V}$  as the similarity metric. Without loss of generality, we use  $Z_{V2T}$  for our similarity metric.

**Scaling for Interpretability** To ensure comparability with metrics like BLEU and ROUGE, we follow a similar approach to CLIP-Score by scaling the embeddings with a weighting factor of 3.5, expanding the score distribution range to [0,100].

## 4 Experiments

To evaluate the effectiveness of SiLVERScore, we conduct multiple experiments to assess the performance compared to back-translation methods.

**Datasets** 1) PHOENIX-14T dataset (Camgoz et al., 2018) is widely recognized as the benchmark dataset for sign language generation (Saunders et al., 2020, 2021; Viegas et al., 2023; Inan et al., 2022). It consists of German Sign Language weather forecast videos segmented into sentences, accompanied by corresponding German transcripts and sign-gloss annotations. 2) CSL-Daily (Zhou et al., 2021). To broaden the domain beyond weather forecasts, we include CSL-Daily, a dataset covering Chinese Sign Language in various daily-life scenarios. This enables us to test the generalizability of SiLVERScore to diverse real-world contexts.

**Translation Model** For the back translation model, we use the multi-stream keypoint attention network proposed by Guan et al., 2024, due to its state-of-the-art performance in sign language translation. This approach minimizes the error propagation caused by inaccuracies in back translation.

**Metrics** We evaluate the quality of back-translated text using both rule-based and embedding-based metrics. For rule-based evaluation, we compute BLEU scores with sacreBLEU (Post, 2018) and ROUGE scores. For embedding-based evaluation, we use BLEURT (specifically BLEURT-20, Pu et al., 2021) and BERTScore (using the bert-base-multilingual-cased model to accommodate the German and Chinese dataset; Zhang\* et al. (2020)). These metrics provide a benchmark for assessing the alignment quality of SiLVERScore in comparison to traditional back-translation evaluation methods.

#### 4.1 Which metric can distinguish between correct and random video-text pairs?

##### 4.1.1 Distribution of Metric Scores

To qualitatively evaluate the performance of different metrics, we analyze the kernel density plots in Figure 2. These plots illustrate the distribution of scores for correctly matched video-text pairs (blue curve) and randomly paired samples (orange curve). SiLVERScore shows a clear separation between the two distributions, with minimal overlap. This indicates its strong ability to distinguish aligned pairs from misaligned ones. In contrast, BLEU-2 exhibits significant overlap, particularly for lower score ranges, suggesting reduced discriminative power for this task. Similarly, the ROUGE shows partial separation but retains overlap between the two distributions. BERTScore and BLEURT show improved separation compared to rule-based metrics but still exhibit some overlap. The sharp distinction and density clustering of scores in the SiLVERScore plot indicate its effectiveness in capturing semantic alignment between video and text representations. Figure 2 focuses on PHOENIX-14T, but we observe similar trends on CSL-Daily.

##### 4.1.2 Quantifying overlap and separability

To complement the qualitative insights from the kernel density plots, we quantify the ability of each metric to distinguish between correctly aligned and randomly paired samples using overlap percentage and ROC AUC (Receiver Operating Characteristic Area Under the Curve). The results are summarized in Table 1. BLEU-4 is omitted for CSL-Daily because consecutive 4-character n-grams in Chinese can lead to sparse counts, producing NaN in the calculation. Since each metric operates on a different scale, we applied Min-Max normalization to scale all metrics to the [0,1] range for a fair comparison.

**Overlap percentage** Overlap percentage measures how much the distributions of scores for correct and random pairs intersect. Lower overlap percentages indicate better discriminative power. SiLVERScore achieves the lowest overlap 6.85% on PHOENIX-14T, and 7.40% on CSL-Daily. BLEU-1 and ROUGE also show single-digit overlaps on CSL-Daily, yet its kernel density plots show that these distributions remain widely dispersed. BERTScore and BLEURT remain competitive with low overlaps, but neither is consistently smaller than SiLVERScore.

**ROC AUC** ROC AUC measures the metric’s ability to distinguish between the two distributions. Higher ROC AUC values indicate better separability, with a maximum value of 1.0. SiLVERScore attains 0.99 AUC for correct vs. random pairs on both datasets, confirming the separation observed in density plots. Overall, the results show that learned embedding-based metrics (SiLVERScore, BERTScore, BLEURT) outperform rule-based metrics in distinguishing between correctly aligned and misaligned video-text pairs.

#### 4.2 Which metric captures semantic distinctions through targeted changes?

Rule-based metrics are inherently sensitive to the exact ordering of words, even when the overall meaning remains unchanged. To demonstrate this sensitivity, we designed an experiment where GPT-4o (Hurst et al., 2024) was used to reorder words in sentences while preserving their meaning. The exact prompt provided to GPT-4o was:

*Reorder the words in the following sentence while keeping the meaning the same: {text} Reordered sentence:*

**Kernel density plot** Figure 3 illustrates how different metrics respond to surface-level changes, specifically word reordering, on PHOENIX-14T. SiLVERScore exhibits the highest score distribution, suggesting its robustness to reordering and its ability to capture semantic content. In contrast, BLEU and ROUGE display sharp peaks and narrower distributions concentrated in the lower score range. This pattern exhibits a clear distinction between rule-based and embedding-based metrics.

**Quantifying overlap and separability** In this experiment, the scores are computed by comparing the ground-truth references with their corresponding hypotheses. While these hypotheses may



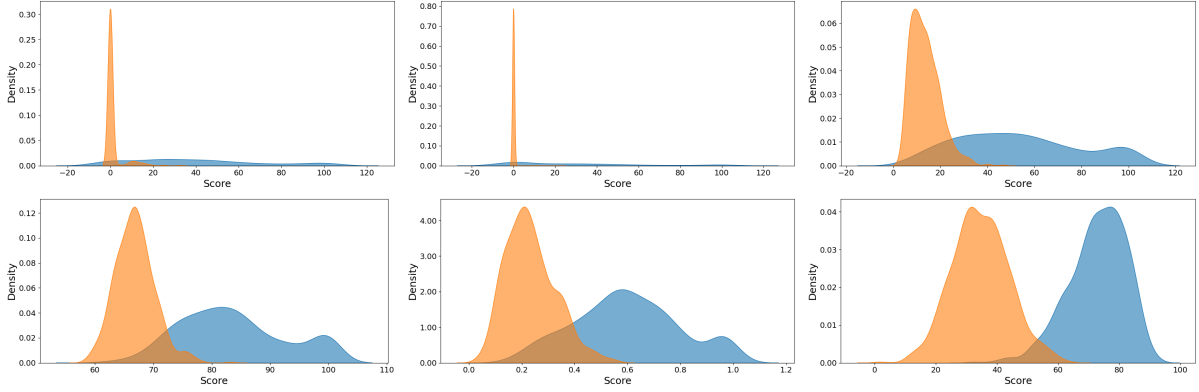


Figure 2: Kernel Density Plots for different metrics. **Top row (left to right, rule-based metrics):** BLEU-2, BLEU-3, ROUGE. **Bottom row (left to right, embedding-based metrics):** BERTScore, BLEURT, SiLVERScore. The blue curve represents the distribution of scores for matching indices (aligned pairs), while the orange curve represents different indices (misaligned pairs). SiLVERScore exhibits a clear separation between the two distributions, indicating a strong ability to distinguish aligned from misaligned pairs. In contrast, BLEU and ROUGE metrics show more overlap, reflecting their sensitivity to surface-level variations.

	Correct vs. Random (§ 4.1)				Original vs. Reordered (§ 4.2)			
	PHOENIX-14T		CSL-Daily		PHOENIX-14T		CSL-Daily	
	Overlap ↓	AUC ↑	Overlap ↓	AUC ↑	Overlap ↑	AUC ↓	Overlap ↑	AUC ↓
BLEU-1	19.78	0.95	6.04	<b>0.99</b>	64.49	0.65	69.28	0.45
BLEU-2	24.30	0.90	5.27	0.98	71.50	0.63	69.60	0.49
BLEU-3	38.63	0.81	23.89	0.88	66.98	0.65	76.43	0.52
BLEU-4	55.45	0.72	-	-	69.47	0.63	83.91	0.53
ROUGE	19.94	0.95	6.12	0.99	67.45	0.67	70.05	0.54
BERTScore	14.17	0.97	9.27	0.98	78.19	0.55	75.18	0.51
BLEURT	21.65	0.95	11.90	0.98	81.31	<b>0.47</b>	70.24	<b>0.39</b>
SiLVERScore	<b>6.85</b>	<b>0.99</b>	<b>7.40</b>	<b>0.99</b>	<b>83.49</b>	0.60	<b>87.84</b>	0.45

Table 1: Comparison of overlap percentages and ROC AUC for various metrics across PHOENIX-14T and CSL-Daily. In “Correct vs. Random” (left columns), lower Overlap and higher AUC reflect better discrimination, and SiLVERScore achieves minimal Overlap (6.85–7.40%) and near-maximal AUC (0.99). In “Original vs. Reordered” (right columns), higher Overlap and lower AUC indicate greater tolerance to meaning-preserving reorderings, where SiLVERScore also achieves the highest Overlap (83.49–87.84%) and lower ROC AUC.

contain errors, they represent the best available approximations of the ground truth. Lower ROC AUC values indicate that the metric maintains its scores despite reordering, reflecting robustness to surface-level variations.

In Table 1, SiLVERScore demonstrates the highest overlap across both datasets (83.49% on PHOENIX-14T, 87.84% on CSL-Daily). This indicates its ability to recognize reordered sentences as semantically equivalent. In contrast, BLEU-1 and BLEU-2 exhibit much lower overlaps, confirming their strict reliance on word order rather than meaning. Moreover, SiLVERScore achieves relatively low ROC AUCs, suggesting it better maintains robustness to reordering.

It is important to note that the original distribu-

tion contains errors, which may affect the Overlap and ROC AUC values for all metrics. This could explain why SiLVERScore’s ROC AUC is slightly higher than those of other metrics.

### 4.3 Which metric can evaluate multimodal and pragmatic aspects more effectively?

#### 4.3.1 Motivation and Setup

Sign languages rely heavily on prosodic markers such as facial expressions, pauses, and intensity to convey meaning. Evaluating the robustness of metrics to prosodic variations is critical, as traditional back-translation-based methods often fail to capture such multimodal cues. We build on the work of Inan et al., 2022, which provided human-annotated token-level prosody intensities for the PHOENIX-

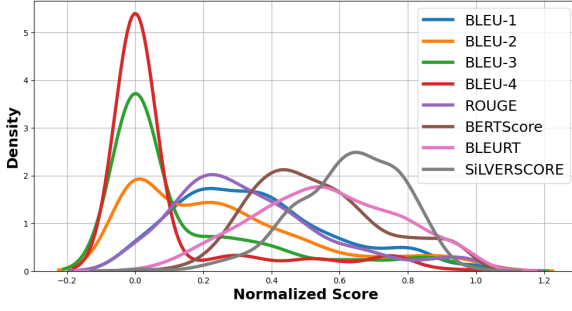


Figure 3: Kernel Density plots comparing the score distributions of different evaluation metrics when applied reordered hypotheses. SiLVERScore, BERTScore, and BLEURT show broader distributions and higher overlap, while rule-based metrics exhibit sharp peaks at lower scores. This indicates their sensitivity to surface-level word order changes.

14T dataset. These annotations classify tokens into three distinct prosodic levels: (i) no intensity: 0, indicating the absence of prosodic markers; (ii) low intensity: 1, reflecting a low degree of intensity markers; and (iii) high intensity: 2, representing high-degree intensity markers.

**Sentence level prosody** We define sentence intensity as the sum of the intensity levels of its tokens,  $I = \sum_{i=1}^n t_i$ , where  $t_i$  is the intensity of token  $i$ . Sentences are categorized into three prosody levels: No Intensity  $I = 0$ , Low Intensity  $1 \leq I \leq 4$ , and High Intensity  $I \geq 5$ .

**Prosody level distribution** The dataset exhibits the following distribution of sentences across these prosody categories: 328 sentences (51.09%) fall under No Intensity, 238 sentences (37.07%) under Low Intensity, and 76 sentences (11.84%) under High Intensity. This distribution indicates that the majority of sentences either lack prosodic markers or exhibit low levels of prosody, while highly expressive sentences are comparatively rare.

#### 4.3.2 Distribution of Scores Across Prosody Categories

To analyze the impact of prosody on evaluation metrics, we categorized sentences based on the sentence-level intensity sums defined earlier. Figure 4 shows the distributions of SiLVERScore, BLEU-1, and ROUGE scores across the categories.

**SiLVERScore Stability** SiLVERScore remains consistent across the three prosody categories, showing minimal variation in median and interquartile range. This demonstrates that SiLVERScore

effectively evaluates semantic alignment without being influenced by prosodic intensity.

**BLEU-1 and ROUGE Sensitivity** BLEU-1 and ROUGE scores decline with increasing prosody intensity, with median scores for High Intensity significantly lower than for No Intensity. This trend indicates that these metrics struggle with prosodically-rich sentences.

**Score Variability** Both BLEU-1 and ROUGE display higher variability in the High Intensity category, suggesting inconsistent performance in evaluating expressive signing.

#### 4.4 Correlation with Prosodic Intensity

As shown in Table 2, traditional back-translation-based metrics (BLEU and ROUGE) exhibit significant negative correlations with prosody intensity (e.g., BLEU-4: -0.200,  $p = 3.31 \times 10^{-7}$ ), reflecting their vulnerability to prosodic variations. This behavior reflects the limitations of traditional metrics, which depend on surface-level text alignment and are vulnerable to information loss during back translation.

Metric	Correlation	p-value
BLEU-1	-0.160	< 0.01
BLEU-2	-0.178	< 0.01
BLEU-3	-0.191	< 0.01
BLEU-4	-0.200	< 0.01
ROUGE	-0.179	< 0.01
BERTScore	-0.144	< 0.01
BLEURT	-0.101	0.01
SiLVERScore	-0.004	0.93

Table 2: Pearson Correlation and p-value of metrics with sentence-level prosody intensity. SiLVERScore demonstrates no significant correlation while other metrics exhibit negative correlations with prosody intensity.

In contrast, SiLVERScore exhibits no significant correlation with prosody intensity (correlation: -0.004,  $p = 0.9277$ ), indicating its robustness to prosodic variations. This robustness suggests SiLVERScore’s ability to evaluate semantic alignment without being influenced by expressive elements.

## 5 The Generalization Problem

While evaluation metrics are expected to generalize across diverse datasets, this remains a significant challenge in sign language processing due to the limited size and diversity of available datasets. As

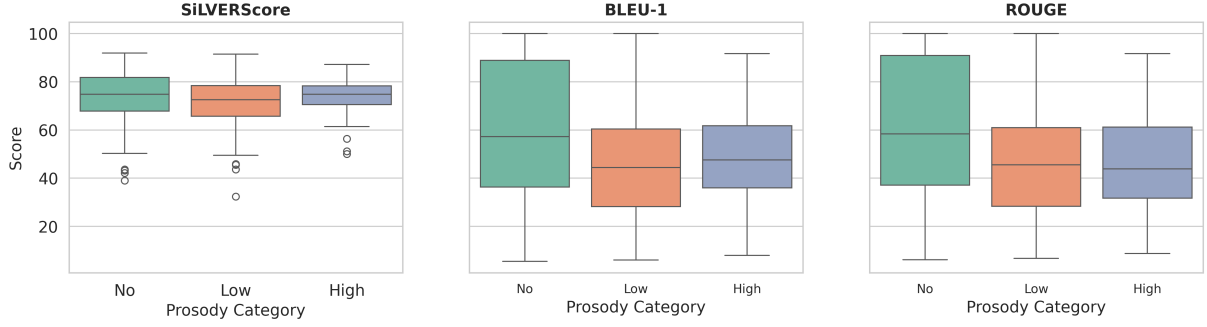


Figure 4: Box plots showing the distribution of SiLVERScore, BLEU-1, and ROUGE scores across three prosody intensity categories (No Intensity, Low Intensity, and High Intensity). While SiLVERScore remains stable across all categories, both BLEU-1 and ROUGE exhibit a noticeable decline in scores as prosody intensity increases. This drop suggests that BLEU-1 and ROUGE are sensitive to prosodically-rich sentences.

highlighted by Jiang et al. (2024), one of the largest sign language datasets, SpreadtheSign, contains only 456,913 examples, which is orders of magnitude smaller than datasets in related domains (e.g., 400M examples for CLIP and 136M for Video-CLIP). In this section, we empirically demonstrate that even SignCLIP, the largest contrastive learning model to date, struggles with generalization at the token level.

## 5.1 Evidence of Limited Generalization

### 5.1.1 Token Level Generalization

We evaluated SignCLIP on ASL Citizen (Desai et al., 2024) and ASL Signs (Chow et al., 2023). The results show that SignCLIP’s generalization capability is limited without fine-tuning.

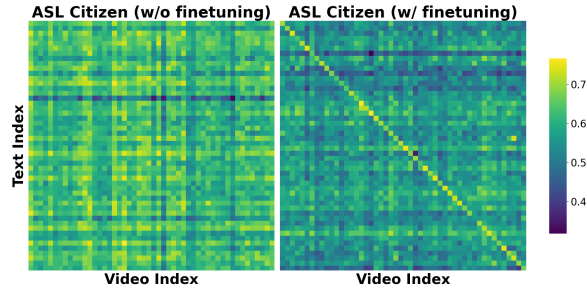


Figure 5: Heatmaps of SignCLIP embeddings cosine similarity scores for two datasets: ASL Citizen (token level) and WMTSLT (sentence level). **Left:** Finetuning increases alignment, as indicated by the clearer diagonal line. **Right:** After finetuning, the model appears to overfit, assigning high similarity scores to many pairs.

Figure 5 illustrates the cosine similarity between video and text embeddings. Ideally, high similarity values should appear along the diagonal, indicating alignment between corresponding video-text pairs. Before fine-tuning, the heatmaps dis-

play low, diffuse similarity scores, indicating poor video-text alignment. Fine-tuning significantly improves alignment, indicating the necessity of dataset-specific adaptation. A similar trend is observed for ASL Signs.

### 5.1.2 Sentence Level Generalization

We evaluated SignCLIP’s sentence-level generalization on the WMTSLT Focus News Corpus (Mathias et al., 2022). Despite fine-tuning, SignCLIP struggles to achieve strong results ( $R@1 = 0.0436$ ).

### 5.1.3 Token Level Language Specific Generalization

To investigate the effect of data size on generalization, we fine-tuned SignCLIP using combined training samples from ASL Signs and SemLex datasets. Despite this, SignCLIP fails to generalize effectively to ASL Citizen ( $R@5 = 0.0005$ ). Even when training on all three datasets, the test set performance on ASL Citizen did not improve significantly. This suggests that dataset-specific characteristics influence performance even when substantial training data is available.

### 5.1.4 Representation Density

Ye et al., 2024 identified a representation density problem, where the semantic visual representations of different sign gestures tend to be closely clustered together, making them hard to distinguish. The proposed contrastive learning strategy, SignCL, encourages the learning of discriminative feature representations. However, applying SignCL to our data yielded limited improvement in retrieval results ( $R@1 = 9.11 \times 10^{-5}$ ), compared to ( $R@1 = 3.04 \times 10^{-5}$ ) with vanilla contrastive learning.

Fine-tuned on	Tested on	R @ 1	R @ 5	R @ 10
<b>Token Level (§ 5.1.1)</b>				
-	Citizen	$1.40 \times 10^{-3}$	$6.10 \times 10^{-3}$	$1.12 \times 10^{-2}$
Citizen	Citizen	$6.39 \times 10^{-2}$	$2.71 \times 10^{-1}$	$4.39 \times 10^{-1}$
<b>Sentence Level (§ 5.1.2)</b>				
WMTSLT	WMTSLT	$3.70 \times 10^{-3}$	$1.75 \times 10^{-2}$	$3.23 \times 10^{-2}$
<b>Token Level Language Specific (§ 5.1.3)</b>				
Signs, SemLex	Citizen	$3.04 \times 10^{-5}$	$5.00 \times 10^{-4}$	$8.00 \times 10^{-4}$
Citizen, Signs, SemLex	Citizen	$4.36 \times 10^{-2}$	$1.76 \times 10^{-1}$	$2.88 \times 10^{-1}$
<b>With SignCL (§ 5.1.4)</b>				
Signs, SemLex	Citizen	$9.11 \times 10^{-5}$	$5.00 \times 10^{-4}$	$9.00 \times 10^{-4}$
<b>With Data Augmentation (§ 5.1.5)</b>				
Signs, SemLex	Citizen	$0.00 \times 10^0$	$2.00 \times 10^{-4}$	$6.00 \times 10^{-4}$
Signs, SemLex	Citizen	$6.07 \times 10^{-5}$	$9.11 \times 10^{-5}$	$3.00 \times 10^{-4}$

Table 3: Text-to-Video Retrieval results and generalization across datasets. Results are shown for different fine-tuning datasets and test datasets.

### 5.1.5 Data Augmentation

Data augmentation is a commonly employed technique to improve model generalization, especially in domains with limited data. To this end, we experimented with several data augmentation strategies including: spatial 2D augmentation, temporal augmentation, and Gaussian noise on keypoints (Jiang et al., 2024). Results show negligible gains (R@1 = 0 with 2D-aug;  $6.07 \times 10^{-5}$  with temporal augmentation), highlighting the limitations of conventional augmentation techniques in enhancing generalization. This suggests that limited dataset diversity and the complexity of visual sign representations cannot be fully addressed through conventional augmentation techniques alone.

## 5.2 How SiLVERScore Addresses Generalization Challenges

Our findings from the experiments suggest the idea that, given current constraints in data availability, tailoring metrics to specific datasets is necessary to create alignment between text and sign.

We proposed a dataset-specific evaluation metric designed to leverage the strengths of embedding-based methods while addressing the constraints of current sign language datasets. By optimizing for specific domains and datasets, we can achieve more reliable evaluations and better alignment with the linguistic and multimodal nature of sign language.

## 6 Conclusion

Through the introduction of SiLVERScore, we demonstrated the empirical strengths of embedding-based methods, including robustness to semantic

variation, prosodic intensity, and a more holistic multimodal evaluation. Our results show that SiLVERScore can overcome limitations of traditional back-translation metrics.

SiLVERScore has the potential to reshape sign language evaluation standards by advancing accessibility for the DHH community and promoting inclusivity in language technologies. Its robustness and semantic sensitivity make it well-suited for broader challenges in multimodal NLP, such as cross-lingual evaluation and integration with video generation models. To support open research and encourage further advancements, we release the code for SiLVERScore’s analysis and computation.

Future efforts should integrate insights from computer graphics, such as improved modeling of spatial relationships and prosody in sign language, to further refine embedding-based methods. Incorporating richer multimodal features, including gesture dynamics and temporal coherence, could enhance the evaluation of expressive and context-dependent signing. Additionally, addressing the scarcity of diverse, large-scale datasets remains critical for improving model generalization.

## Acknowledgments

This research was supported in part by the U.S. National Science Foundation under Award No. 2418664. We thank Asteria Kaeberlein and Katherine Atwell for their helpful feedback.

## References

Danielle Bragg, Oscar Koller, Mary Bellard, Larwan Berke, Patrick Boudreault, Annelies Braffort, Naomi



- Caselli, Matt Huenerfauth, Hernisa Kacorri, Tessa Verhoef, et al. 2019. Sign language recognition, generation, and translation: An interdisciplinary perspective. In *Proceedings of the 21st international ACM SIGACCESS conference on computers and accessibility*, pages 16–31.
- Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. 2018. Neural sign language translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yiting Cheng, Fangyun Wei, Jianmin Bao, Dong Chen, and Wenqiang Zhang. 2023. [CiCo: Domain-Aware Sign Language Retrieval via Cross-Lingual Contrastive Learning](#). In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19016–19026, Los Alamitos, CA, USA. IEEE Computer Society.
- Ashley Chow, Glenn Cameron, Mark Sherwood, Phil Culliton, Sam Sepah, Sohler Dane, and Thad Starner. 2023. Google - isolated sign language recognition. <https://kaggle.com/competitions/asl-signs>. Kaggle.
- Aashaka Desai, Lauren Berger, Fyodor O. Minakov, Vanessa Milan, Chinmay Singh, Kriston Pumpfrey, Richard E. Ladner, Hal Daumé, Alex X. Lu, Naomi Caselli, and Danielle Bragg. 2024. Asl citizen: a community-sourced dataset for advancing isolated sign language recognition. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA. Curran Associates Inc.
- Lu Dong, Lipisha Chaudhary, Fei Xu, Xiao Wang, Mason Lary, and Ifeoma Nwogu. 2024. [Signavatar: Sign language 3d motion reconstruction and generation](#). In *2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–10.
- Mo Guan, Yan Wang, Guangkun Ma, Jiarui Liu, and Mingzu Sun. 2024. [Multi-stream keypoint attention network for sign language recognition and translation](#). *ArXiv*, abs/2405.05672.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. [CLIPScore: A reference-free evaluation metric for image captioning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, page 6629–6640, Red Hook, NY, USA. Curran Associates Inc.
- Wencan Huang, Wenwen Pan, Zhou Zhao, and Qi Tian. 2021. Towards fast and high-quality sign language production. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 3172–3181.
- Matt Huenerfauth, Liming Zhao, Erdan Gu, and Jan Allbeck. 2008. Evaluation of american sign language generation by native asl signers. *ACM Transactions on Accessible Computing (TACCESS)*, 1(1):1–27.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Eui Jun Hwang, Jung Ho Kim, Suk Min Cho, and Jong C Park. 2022. Non-autoregressive sign language production via knowledge distillation. *arXiv preprint arXiv:2208.06183*.
- Eui Jun Hwang, Huije Lee, and Jong C. Park. 2024. [A gloss-free sign language production with discrete representation](#). In *2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–6.
- Mert Inan, Yang Zhong, Sabit Hassan, Lorna Quandt, and Malihe Alikhani. 2022. [Modeling intensification for sign language generation: A computational approach](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2897–2911, Dublin, Ireland. Association for Computational Linguistics.
- Zifan Jiang, Gerard Sant, Amit Moryossef, Mathias Müller, Rico Sennrich, and Sarah Ebling. 2024. [Sign-CLIP: Connecting text and sign language by contrastive learning](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9171–9193, Miami, Florida, USA. Association for Computational Linguistics.
- Jung-Ho Kim, Mathew Huerta-Enochian, Changyong Ko, and Du Hui Lee. 2024. SignBLEU: Automatic evaluation of multi-channel sign language translation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, Turin, Italy.
- Harold W Kuhn. 1955. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97.
- Scott K. Liddell. 2003. *Grammar, Gesture, and Meaning in American Sign Language*. Cambridge University Press.
- Li Liu, Lufei Gao, Wentao Lei, Fengji Ma, Xiaotian Lin, and Jinting Wang. 2023. A survey on deep multi-modal learning for body language recognition and generation. *arXiv preprint arXiv:2308.08849*.

- Müller Mathias, Ebling Sarah, Camgöz Necati Cihan, Jiang Zifan, Battisti Alessia, Moryossef Amit, Rios Annette, Bowden Richard, and Wong Ryan. 2022. [Wmt-slt focusnews: Training data for the wmt shared task on sign language translation](#).
- Amit Moryossef, Rotem Zilberman, and Ohad Langer. 2024. signwriting-evaluation: Effective sign language evaluation via signwriting. *arXiv preprint arXiv:2410.13668*.
- B Natarajan, E Rajalakshmi, R Elakkiya, Ketan Kotecha, Ajith Abraham, Lubna Abdelkareim Gabralla, and V Subramaniaswamy. 2022. Development of an end-to-end deep learning framework for sign language recognition, translation, and video generation. *IEEE Access*, 10:104358–104374.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Amy Pu, Hyung Won Chung, Ankur P Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for mt. In *Proceedings of EMNLP*.
- Wendy Sandler. 2012. [The phonological organization of sign languages](#). *Language and Linguistics Compass*, 6(3):162–182. Epub 2012 Mar 2.
- Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. 2020. [Progressive transformers for end-to-end sign language production](#). In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI*, page 687–705, Berlin, Heidelberg. Springer-Verlag.
- Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. 2021. [Continuous 3d multi-channel sign language production via progressive transformers and mixture density networks](#). *Int. J. Comput. Vision*, 129(7):2113–2135.
- Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. In *Proceedings of ACL*.
- Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. 2019. Fvd: A new metric for video generation.
- Gül Varol, Liliane Momeni, Samuel Albanie, Triantafyllos Afouras, and Andrew Zisserman. 2021. [Read and attend: Temporal localisation in sign language videos](#). In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16852–16861.
- Carla Viegas, Mert Inan, Lorna Quandt, and Malihe Alikhani. 2023. [Including facial expressions in contextual embeddings for sign language generation](#). In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (\*SEM 2023)*, pages 1–10, Toronto, Canada. Association for Computational Linguistics.
- Pan Xie, Taiying Peng, Yao Du, and Qipeng Zhang. 2024. [Sign Language Production with Latent Motion Transformer](#). In *2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 3012–3022, Los Alamitos, CA, USA. IEEE Computer Society.
- Jinhui Ye, Xing Wang, Wenxiang Jiao, Junwei Liang, and Hui Xiong. 2024. Improving gloss-free sign language translation by reducing representation density.
- Kayo Yin, Amit Moryossef, Julie Hochgesang, Yoav Goldberg, and Malihe Alikhani. 2021. [Including signed languages in natural language processing](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7347–7360, Online. Association for Computational Linguistics.
- Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee. 2020. [Speech gesture generation from the trimodal context of text, audio, and speaker identity](#). *ACM Transactions on Graphics (TOG)*, 39:1 – 16.
- Tianyi Zhang\*, Varsha Kishore\*, Felix Wu\*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.
- Hao Zhou, Wengang Zhou, Weizhen Qi, Junfu Pu, and Houqiang Li. 2021. Improving sign language translation with monolingual data by sign back-translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1316–1325.