

Beyond a Single Reference: Training and Evaluation with Paraphrases in Sign Language Translation

Anonymous ACL submission

Abstract

Most Sign Language Translation (SLT) corpora pair each signed utterance with a single written-language reference, despite the highly non-isomorphic relationship between sign and spoken languages, where multiple translations can be equally valid. This limitation constrains both model training and evaluation, particularly for n-gram-based metrics such as BLEU.

In this work, we investigate the use of Large Language Models to automatically generate paraphrased variants of written-language translations as synthetic alternative references for SLT. First, we compare multiple paraphrasing strategies and models using an adapted ParaScore metric. Second, we study the impact of paraphrases on both training and evaluation of the pose-based T5 model on the YouTubeASL and How2Sign datasets.

Our results show that naively incorporating paraphrases during training does not improve translation performance and can even be detrimental. In contrast, using paraphrases during evaluation leads to higher automatic scores and better alignment with human judgments. To formalize this observation, we introduce BLEU_{para}, an extension of BLEU that evaluates translations against multiple paraphrased references. Human evaluation confirms that BLEU_{para} correlates stronger with perceived translation quality. We release all generated paraphrases, generation and evaluation code to support reproducible and more reliable evaluation of SLT systems.

1 Introduction

Automatic Sign Language Translation (SLT) has recently made rapid progress due to advances in Vision Language Models and the availability of large-scale video-text corpora. Despite these developments, a persistent limitation in current SLT datasets is the scarcity of translation variability: each sign language utterance is typically paired

with a single weakly aligned written-language reference translation. This stands in contrast to spoken language machine translation, where corpora frequently provide multiple reference translations or explicitly annotate translation alternatives, capturing the natural diversity of phrasing, register, and information structure that exists in human language. Such variation is crucial for both training and evaluation, where metrics such as BLEU are known to be sensitive to the availability of reference translations due to their n-gram matching nature.

Sign languages pose an even stronger requirement for translation diversity. The mapping between visual expression and written language is highly non-isomorphic: a single sign language sentence can be translated into several equally valid spoken language realizations, differing in word order, degree of explicitness, and syntactic structure. Relying on a single reference translation therefore risks over-constraining model training and underestimating translation quality during evaluation.

In this work, we explore whether Large Language Models (LLMs) can be used to systematically enrich SLT corpora with automatically generated paraphrases of written language translations. Our goal is twofold: (1) to study whether synthetic generation of translation variants (i.e. paraphrasing) can improve the robustness of sequence-to-sequence SLT models such as T5, and (2) to examine how these paraphrases affect evaluation, comparing BLEU scores against human evaluation.

By viewing paraphrases as alternative translations, analogous to multi-reference machine translation setups, we aim to bridge principles from classical natural language processing (NLP) with the specific challenges of SLT.

2 Related work

Sign Language Translation Recent progress in SLTn have been driven by increasingly powerful

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

multimodal and transformer-based architectures. A commonly used baseline architecture consists of a keypoint extractor (e.g. MMPose (Sengupta et al., 2020) or MediaPipe (Lugaresi et al., 2019)), followed by transformer (Vaswani et al., 2017). This approach became an effective way to perform experiments and provide benchmark results for new datasets (Uthus et al., 2023; Camgöz et al., 2021; Tanzer and Zhang, 2025; Tanzer, 2025; Alhejab et al., 2025). Recent state-of-the-art methods extend this approach by adding additional modules. GloFE (Lin et al., 2023) employs visual backbone based on graph convolutional networks (GCN) to more effectively process pose features, followed by a transformer encoder–decoder architecture that learns robust gloss-level representations through contrastive learning. Uni-Sign (Li et al., 2025) similarly uses GCN for pose feature encoding and further incorporates a fusion module that integrates RGB information in cases where keypoint extraction is failing. SignBERT+ (Hu et al., 2023) extends BERT-style pretraining (Devlin et al., 2019) to sign language data, learning contextualized representations that capture temporal dependencies within signing sequences and can be transferred to downstream tasks such as isolated and continuous sign language recognition. LITFIC (Jang et al., 2025) leverages a LLaMA-based (Touvron et al., 2023) large language model that incorporates contextual and meta-information, including pseudo-gloss cues, to improve the translation quality. By integrating these diverse sources of information, the model can better capture discourse-level context and disambiguate signs.

Paraphrasing Paraphrasing - generating semantically equivalent expressions with varied surface forms - has been extensively studied for both its theoretical properties and practical applications in NLP. Bhagat & Hovy (Bhagat and Hovy, 2013) established that quality paraphrases balance semantic preservation against lexical and syntactic diversity, a tension later formalized in evaluation metrics like ParaScore (Shen et al., 2022), which combines BERTScore-based (Zhang et al., 2020) semantic similarity with explicit diversity modeling. Neural paraphrase generation has evolved from seq2seq models to transformer-based approaches fine-tuned on large-scale datasets like ParaNMT-50M (Iyyer et al., 2018). In machine translation, paraphrasing serves as a key data augmentation strategy, particularly for low-resource settings where back-

translation (Sennrich et al., 2016) and lexical substitution have yielded substantial BLEU improvements. For SLT specifically, paraphrasing is particularly relevant where due to the structural divergence between signed and written languages means a single signed utterance often admits multiple valid translations, yet datasets typically provide only one translation.

Data In this research, we focus on American Sign Language, namely by using two widespread datasets - YoutubeASL (Uthus et al., 2023) and How2Sign (Duarte et al., 2021). YouTubeASL is a large-scale dataset of American Sign Language (ASL) collected from publicly available YouTube videos. It contains over 610,000 video samples (equivalent to about 1,000 hours) and features 2,500 unique signers from diverse backgrounds. Videos cover various domains, such as education, personal vlogs, and conversational ASL, making the dataset representative of a wide range of signing contexts.

How2Sign is a smaller but well-annotated dataset, originally developed to study sign translation under multiple angles/views. It contains 35,000 signing samples, amounting to approximately 80 hours of ASL content, produced by roughly 11 signers in a controlled environment with stable lighting and fixed backgrounds. Although the dataset’s primary focus is translating from spoken English to ASL, it is widely adopted for other tasks in the sign language processing community, including sign language translation. One challenge with using How2Sign lies in the presence of translationese effects, a phenomenon where translated text can sound artificial or lack natural linguistic flow (Graham et al., 2020). Despite this limitation, How2Sign serves as an important benchmark for evaluating model performance. Its controlled nature provides consistency that complements the more variable, in-the-wild data from YouTubeASL.

3 Methods

3.1 Sign Language Translation

For our experiments, we adopt an existing SLT framework that uses a modified T5 encoder–decoder transformer to process pose-based inputs and includes comprehensive data preprocessing, such as keypoint extraction, normalization, interpolation, and augmentation (Zelezny et al., 2025). We adopt this existing setup because it

provides already ready-to-use codebase and preprocessed YouTubeASL dataset (Zelezny et al., 2024), including extracted keypoints and corresponding annotations. This allows us to focus on evaluating our novel contributions without re-implementing core components such as data handling and model configuration. The publicly available repository accompanying that work includes all preprocessing pipelines, model training routines, and evaluation scripts, ensuring reproducibility and consistent comparison across experiments.

3.2 Paraphrasing and ParaScore

Paraphrasing plays a vital role in modern natural language processing and machine translation, making the quality assessment of paraphrases essential for these applications. In the context of SLT, where the mapping between visual and written modalities admits multiple valid realizations, paraphrases can serve as synthetic valid alternative references. However, not all automatically generated paraphrases will be equally useful: a paraphrase that diverges too far semantically fails to represent the source meaning, while one that is too lexically similar to the original provides little additional signal. A successful paraphrase must balance two critical criteria: semantic similarity and lexical divergence. To evaluate this balance, we employ an adapted version of the ParaScore metric (Shen et al., 2022), which effectively combines semantic preservation with linguistic diversity through the use of both language models and traditional linguistic measures.

Unlike most existing metrics that struggle to adequately balance semantic similarity and lexical divergence, this approach implements a max function for the semantic similarity component between an input sequence and its candidate rephrase-ment. The semantic similarity component utilizes BERTScore (Zhang et al., 2020), which employs the following process: it first computes contextual embeddings for each token in both sentences, then calculates cosine similarities between tokens in the two sentences, and finally employs a matching strategy to compute precision, recall, and F1 score.

For measuring lexical divergence, we implement the Levenshtein distance (Levenshtein, 1965), which calculates the minimum number of single-character edits (insertions, deletions, or substitutions) required to transform one word into another. This distance is normalized by dividing by the length of the longer string, resulting in a value be-

tween 0 and 1, where 0 indicates identical strings and 1 represents completely different strings.

The metric incorporates two hyperparameters: γ and ω . γ , set to a default value of 0.35, serves two purposes: it caps the maximum divergence score and introduces non-linearity to the scaling of the divergence score. This means that edit distances above 35% of the longer string’s length are treated equivalently. ω , with a default value of 0.5, is used in the final ParaScore calculation to weight the divergence. Formally, given an input x and a candidate paraphrase \hat{x} , the metric is defined as:

$$\text{ParaScore}(x, \hat{x}) = \frac{\text{BERTScore}(x, \hat{x}) + \omega \cdot \min(\text{NLD}(x, \hat{x}), \gamma)}{1 + (\omega \cdot \gamma)} \quad (1)$$

where NLD is the Normalized Levenshtein Distance. With these parameters, the similarity score can contribute up to 1 to the final score, while the divergence score can contribute up to 0.175 (i.e., 0.5×0.35). Although this creates a maximum possible score of 1.175, we normalize the final score by dividing by this maximum (represented by the denominator) to achieve a standardized scale.

4 Comparison of Paraphrasing Methods

In this section we compare paraphrasing models under a single, controlled generation pipeline. For each reference English sentence x in our SLT corpora, we ask a candidate LLM to produce $K=5$ meaning-preserving rewrites $\{\hat{x}_k\}_{k=1}^K$ using the same system instruction across models: You are a helpful assistant that rephrases a given sentence in K ways, each on its own line. Try to be semantically consistent and output nothing else than these sentences. Sentence: $\langle x \rangle$ Paraphrases: . To keep outputs comparable, we use consistent decoding (sampling; temperature = 0.7, top- $p=0.95$). We then normalize generations by stripping list markers/bullets and occasional boilerplate text (e.g., “Here are ...”); sentences outside 4–30 words are skipped, and any generation that does not yield exactly K lines is treated as missing. The resulting per-model paraphrase sets are scored with ParaScore and used for the analyses in Figures 1, 2, 3 and the subsequent experiments.

Through careful manual evaluation of paraphrases at various percentile levels (0.25, 0.50, 0.75), we established a quality threshold of 0.7. This threshold serves as a reliable indicator for dis-

280 distinguishing high-quality paraphrases from lower- 330
281 quality ones. 331

282 For paraphrasing ASL translations, we use a 332
283 wide range of current Large Language Models. 333
284 Their results are compared in Figure 1. 334

285 Because SLT is currently done mostly at 335
286 sentence-level, the given sentence is usually a 336
287 few seconds clip, extracted from a longer video, 337
288 lacking context, which was proven particularly 338
289 problematic as Sign Languages are very context- 339
290 dependent (Jang et al., 2025). 340

291 As a separate experiment, we augment the para- 341
292 phrasing prompt with lightweight *video-level tex-* 342
293 *tual context*. For each target sentence, we add 343
294 the preceding reference sentence(s) from the same 344
295 video clip, and reset the context when moving to a 345
296 new video. This context is provided as part of the 346
297 instruction while keeping the target sentence un- 347
298 changed, allowing us to isolate whether discourse 348
299 information (co-reference, ellipsis, topic conti- 349
300 nuity) improves meaning-preserving paraphrases 350
301 compared to sentence-only prompting. However, 351
302 in our experimental setting the ParaScore is worse 352
303 for the paraphrases generated with context (see 353
304 Figure 2). This is something we are planning to 354
305 explore further in our future work. 355

306 Another separate experimental setting tackles 356
307 the use of prompting to gain multiple paraphrases 357
308 of the same reference translation sentence. We can 358
309 either use **sequential** (ask for n paraphrases at once 359
310 in a single prompt) or **iterative** prompting (ask 360
311 for a paraphrase n times). Sequential prompting 361
312 proved to generate paraphrases with a little higher 362
313 average ParaScore (see Figure 3) and for a lower 363
314 computation cost. 364

315 In our comparison, we applied our adapted 365
316 ParaScore metric to each dataset, analyzing both 366
317 individual paraphrases and the average scores for 367
318 each caption. The analysis revealed that GPT-4o- 368
319 mini paraphrasing achieved the highest average 369
320 ParaScore; therefore, its generated paraphrase sets 370
321 for YouTubeASL and How2Sign are used in the 371
322 subsequent training and evaluation experiments. 372
323 Notably, we observed that adding context to both 373
324 sequential and iterative prompting in LLaMA led 374
325 to significant performance degradation. 375

326 5 Paraphrases for Training and 376 327 Evaluation 377

328 In this section, we investigate the impact of para- 378
329 phrases on SLT with two main objectives: (1) To 379

examine whether incorporating paraphrased tar- 330
get sentences during training improves SLT per- 331
formance. (2) To assess whether using paraphrases 332
can lead to improved automatic evaluation quality. 333

334 Different paraphrase-based training and evalua- 335
336 tion strategies are compared in a controlled experi- 336
337 mental setting. Evaluation using paraphrased refer- 337
338 ences is particularly interesting, as higher align- 338
339 ment with human judgment would indicate their 339
340 potential usefulness; the analysis of human judg- 340
ments is discussed in Section 6.

341 5.1 Training Setup 341

342 We follow the default training configuration of the 342
343 adopted framework (see 3.1), using the same hy- 343
344 perparameter values and preprocessing settings. 344
345 Experiments are conducted on the pose-based 345
346 YouTubeASL dataset with pre-extracted keypoints, 346
347 adopting the provided 90:10 train-val split. For 347
348 the How2Sign dataset, we use the original dataset 348
349 split. Input pose sequences are normalized using 349
350 SignSpace (Zelezny et al., 2025) normalization, 350
351 with missing keypoints filled using a fixed value 351
352 of -10 . Models are pretrained on YouTubeASL 352
353 dataset with an effective batch size of 256 and 353
354 a learning rate of 0.0004, employing a warm-up 354
355 phase of 5,000 steps. The models are trained until 355
356 convergence, i. e. about 400,000 iteration steps. 356
357 We then finetune the models on How2Sign dataset 357
358 using the same setting until convergence, i. e. about 358
359 10,000 steps. No data augmentation techniques are 359
360 applied during nor of the training stages. 360

361 For evaluation, we report BLEU scores com- 361
362 puted using the sacrebleu library v2.4.3¹, 362
363 ROUGE-L scores from evaluate library², and 363
364 BLEURT scores using the BLEURT-20 checkpoint 364
365 from the official repository³. 365

366 5.2 Experiments 366

367 First, the model is pretrained on the YouTube ASL 367
368 dataset without paraphrases. It is then fine-tuned 368
369 on the How2Sign dataset using different types of 369
370 paraphrase-based training strategies. The result- 370
371 ing model checkpoints are then evaluated using 371
372 the How2Sign test set. To ensure consistency of 372
373 results, each experiment is run three times with 373
374 different random seeds, and all reported results cor- 374
375 respond to the average performance across these 375

¹<https://pypi.org/project/sacrebleu/2.4.3/>

²<https://github.com/huggingface/evaluate>

³<https://github.com/google-research/bleurt>

ParaScore Distribution Curves by Model

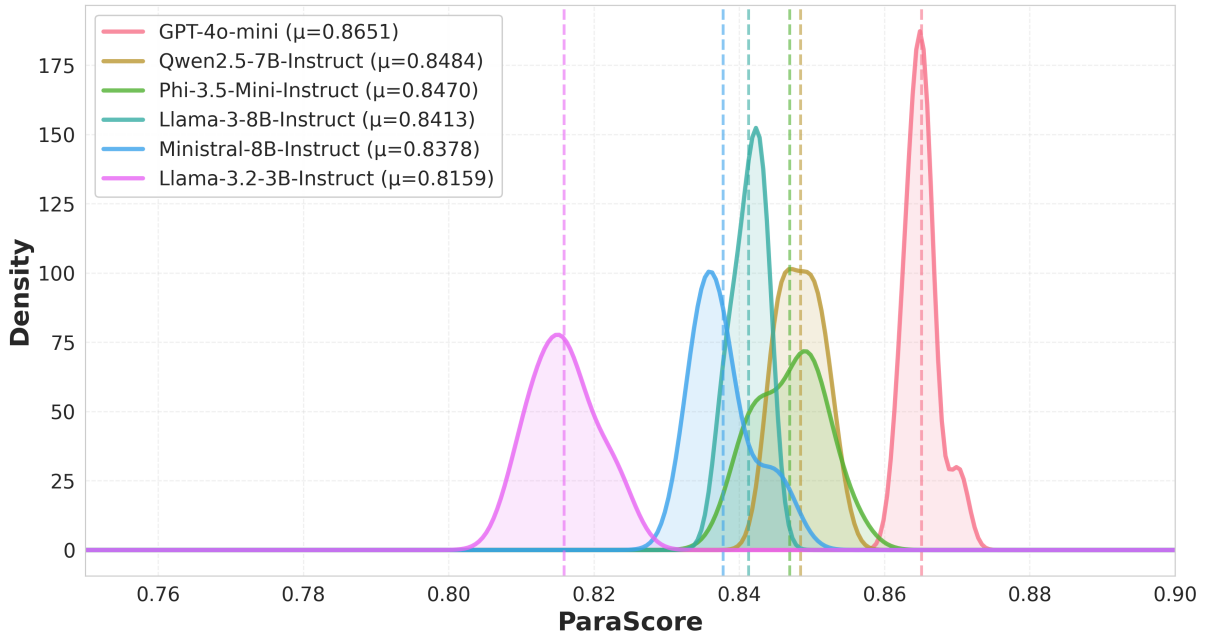


Figure 1: **ParaScore distributions for paraphrases generated by different LLMs.** Kernel density estimates over all generated paraphrases (same prompt and decoding for all models). Dashed vertical lines mark per-model means (μ ; shown in the legend). Higher scores indicate paraphrases that better preserve meaning while avoiding near-copies, enabling a direct quality comparison across paraphrasing models.

376 runs. Complete results are provided in the supple-
 377 mentary material.

378 We consider three training configurations. First,
 379 we train a baseline model using only the canonical
 380 target sentence without paraphrases. Second, we
 381 train a model where the target sentence is randomly
 382 sampled from up to six alternatives (one canonical
 383 reference and five paraphrases) for each training in-
 384 stance. Third, we compute the training loss for all
 385 available paraphrases and backpropagate the gra-
 386 dients from the paraphrase yielding the minimum
 387 loss for a given instance.

388 In the paraphrase-based evaluation setting, we
 389 calculate metrics for all available paraphrases of
 390 each test instance and select the highest-scoring
 391 paraphrase as the final result. As shown in Ta-
 392 ble 1, models trained without paraphrases outper-
 393 form those trained with paraphrased targets. This
 394 suggests that exposure to multiple paraphrased tar-
 395 gets during training may introduce ambiguity that
 396 negatively impacts the performance. As expected,
 397 evaluation scores obtained using paraphrased refer-
 398 ences are generally higher. However, the relation-
 399 ship between these automatic metrics and human
 400 judgment is analyzed in the following section.

Par. mode	BLEU-4	BLEURT	Rouge-L
Evaluation without paraphrases			
No par.	7.46	0.43	0.29
Random	6.17	0.43	0.26
Min loss	6.49	0.43	0.27
Evaluation with paraphrases			
No par.	8.86	0.44	0.31
Random	7.95	0.45	0.30
Min loss	8.00	0.44	0.30

Table 1: Performance comparison of different paraphrase-based training strategies on the How2Sign test set using BLEU-4, BLEURT, and ROUGE-L metrics. The checkpoints were evaluated both without and with paraphrases, where metrics are computed for all available paraphrases of each test instance and the highest-scoring paraphrase is selected. No par. denotes training with only canonical targets, Random indicates random sampling from available paraphrases during training, and Min loss selects the paraphrase yielding the minimum training loss per instance. For the Min loss configuration, a paraphrased target was selected in approximately 60% of training instances.

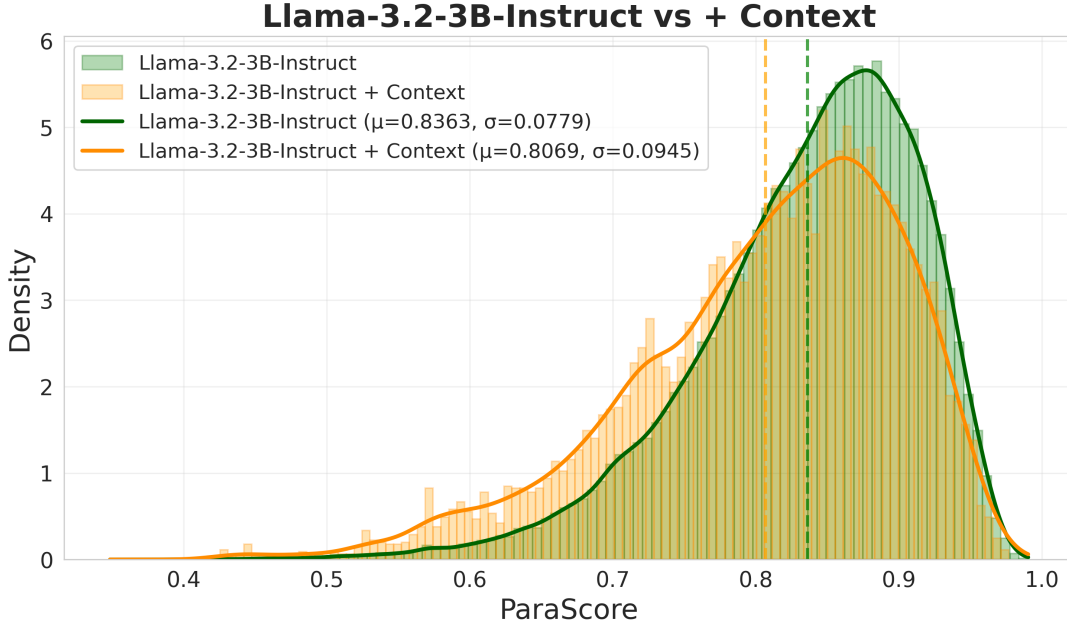


Figure 2: **Effect of adding video-level context to the paraphrasing prompt.** ParaScore distributions for Llama-3.2-3B-Instruct with sentence-only prompting versus prompting augmented with short preceding context from the same video clip. Dashed lines denote means (μ); providing context shifts the distribution left and increases variance, suggesting that extra discourse information can encourage looser rewrites that more often drift from the original meaning under our ParaScore criterion.

6 Human Evaluation

To address prior findings (Jang et al., 2025), that show concerning limitations of the standardized metrics in SLT, usually focused on WER and translation length, we propose the $\text{BLEU}_{\text{para}}$ metric, enhancing the standard BLEU4 (Papineni et al., 2002) metric with 5 paraphrases of each reference sentence in the dataset, fixing any word order or synonym issues. To enable its use, we publish the LLM-generated paraphrased versions of How2Sign (Duarte et al., 2021) and YoutubeASL (Uthus et al., 2023) alongside this paper.

To prove that this metric improves SLT evaluation over its non-paraphrased counterpart, we have conducted a Human Evaluation experiment, using the exact same protocol as (Jang et al., 2025). In addition to reporting correlations over all items, we explicitly probe *extreme* cases where reference-matching metrics are typically the most brittle: we select the 48 examples whose canonical-reference BLEU is either very low (< 5) or very high (> 15). This subset captures both clear failures and near-verbatim matches, and highlights whether $\text{BLEU}_{\text{para}}$ better tracks human judgments in outlier translations.

The results in Table 2 show that human percep-

Metric	Pearson r	Spearman ρ	ρ (extremes)
BLEU	0.688	0.657	0.578
$\text{BLEU}_{\text{para}}$	0.697	0.685	0.659

Table 2: Correlation between mean human ratings (0–5; 70 sentences, 6 annotators; mean = 1.45 ± 1.14) and automatic metrics. “Extremes” uses translations with $\text{BLEU} < 5$ or > 15 .

tion of model capabilities correlates better with $\text{BLEU}_{\text{para}}$, especially in outlier cases (48 translations with $\text{BLEU} < 5$ or > 15).

7 Discussion

Our comparison of paraphrasing methods reveals substantial variability among individual LLMs in their ability to generate high-quality paraphrases for synthetic translation data. Among the evaluated models, GPT-4o outperforms the other methods by a considerable margin. We assume the trend to be improving, with the expectation that commercial models will be slightly more capable than publicly available models. It is therefore likely that the results achieved by GPT-4o in our experiments may already be surpassed by newer versions of commercial models, such as GPT-5 or Gemini 3. Nevertheless, because these experiments are finan-

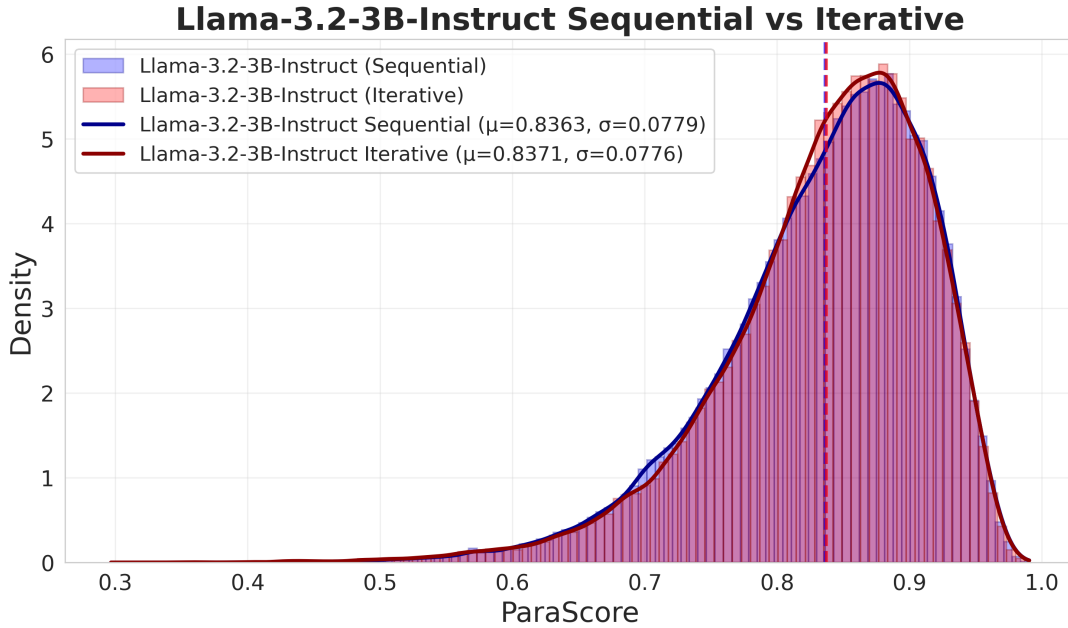


Figure 3: **Sequential vs. iterative paraphrase generation.** ParaScore distributions for Llama-3.2-3B-Instruct when producing five paraphrases in a single call (*sequential*) versus generating them one-by-one with feedback from earlier outputs (*iterative*). Histograms and density curves largely overlap, and the mean/variance (dashed lines; μ, σ in legend) are nearly identical, indicating that iterative prompting does not materially change paraphrase quality in our setting.

445 cially demanding and the primary goal of this paper
 446 is not to achieve state-of-the-art results, but rather
 447 to explore a general concept, we leave additional
 448 experiments for future research.

449 Incorporating paraphrases directly as data aug-
 450 mentation during model training does not yield
 451 measurable performance gains, suggesting that
 452 naive output-level augmentation is ineffective.
 453 While introducing multiple paraphrase variants into
 454 the loss computation could potentially improve
 455 learning, this approach significantly increases train-
 456 ing time in the current setup. More selective strate-
 457 gies, such as hard example mining, may therefore
 458 be a more promising direction for future research.

459 Prior work has demonstrated that current stan-
 460 dardized metrics, such as BLEU and ROUGE, rely
 461 heavily on word-level matching strategies, which
 462 correlate with human perception only weakly. To
 463 address this limitation, we propose a novel metric
 464 $BLEU_{para}$. Human evaluation confirms that the cor-
 465 relation is improved by employing this new metric
 466 instead of the standard BLEU metric. Additionally,
 467 $BLEU_{para}$ partially alleviates issues related to syn-
 468 onyms and word order. Its main disadvantage is the
 469 necessity to have paraphrases for each test sample.
 470 To reduce this burden and support reproducibility,
 471 we publish all paraphrases, including the codes for

their creation, in our GitHub repository .

8 Conclusion

472
 473
 474 This paper experiments with the use of paraphrases
 475 for training and evaluation. We first compare the
 476 capabilities of different large language models in
 477 the task of paraphrase creation. Then we utilize
 478 these paraphrases for both training and evaluation.
 479 While the training with paraphrases, using basic
 480 ideas for its incorporation, does not bring any mea-
 481 surable improvements, we would like to explore
 482 more complex protocols in our future research.

483 On the other hand, the evaluation using para-
 484 phrases provides results better correlated with hu-
 485 man perception of model capabilities. To reflect
 486 the usage of paraphrases during the evaluation step,
 487 we propose a novel metric $BLEU_{para}$. We argue
 488 that this metric is more suitable for the evaluation
 489 of methods for SLT than the standard metrics like
 490 BLEU or ROUGE. All the codes and paraphrases
 491 are publicly available in our GitHub repository .

492 In our future work, we would like to explore
 493 different protocols for paraphrase creation more in
 494 depth, especially the utilization of context. Addi-
 495 tionally, we would like to test if the proposed met-
 496 ric $BLEU_{para}$ generalizes well also for other sign
 497 languages other than American Sign Language.

498 Limitations

499 The main limitation of this research is scope and
500 reproducibility. Our experiments are confined to
501 a small set of datasets and one translation setting,
502 so the findings may not transfer to other sign lan-
503 guages, domains, or annotation styles. Results are
504 also sensitive to the choice of paraphrasing model,
505 decoding settings, and prompts; in particular, us-
506 ing proprietary LLMs and stochastic sampling can
507 make exact reproduction difficult and may change
508 over time as providers update models. Finally, com-
509 pute and cost constraints limited the number of
510 paraphrases per reference, the breadth of hyperpa-
511 rameter sweeps, and the scale of human evaluation,
512 so some effects (especially small gains between
513 strong models) may be underpowered and should
514 be revalidated on larger samples.

515 Our work is limited by the technical properties
516 of sign language video data and its alignment with
517 translations. Variations in resolution, frame rate,
518 motion blur, camera framing, lighting, occlusions,
519 and background contrast can reduce the visibility
520 of important manual and non-manual cues. Video
521 segments are often temporally incomplete or mis-
522 aligned with their translations, lacking natural on-
523 set and offset phases, which introduces noise and
524 complicates training and evaluation. These factors
525 affect both the reliability of model learning and the
526 interpretability of automatic metrics, including ex-
527 periments with paraphrased targets and evaluation
528 references.

529 References

530 Ali Alhejab, Tomas Zelezny, Lamy Alkanhal, Ivan
531 Gruber, Yazeed Alharbi, Jakub Straka, Vaclav Ja-
532 vorek, Marek Hruz, Badriah Alkalifah, and Ahmed
533 Ali. 2025. Saudi sign language translation using
534 t5. In *International Conference on Speech and Com-
535 puter*, pages 331–343.

536 Rahul Bhagat and Eduard Hovy. 2013. [Squibs: What is
537 a paraphrase?](#) *Computational Linguistics*, 39(3):463–
538 472.

539 Necati Cihan Camgöz, Ben Saunders, Guillaume Ro-
540 chette, Marco Giovanelli, Giacomo Inches, Robin
541 Nachtrab-Ribback, and Richard Bowden. 2021. [Con-
542 tent4all open research sign language translation
543 datasets](#). In *2021 16th IEEE International Confer-
544 ence on Automatic Face and Gesture Recognition
545 (FG 2021)*, pages 1–5.

546 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
547 Kristina Toutanova. 2019. Bert: Pre-training of deep

bidirectional transformers for language understand- 548
ing. In *Proceedings of the 2019 conference of the 549
North American chapter of the association for com- 550
putational linguistics: human language technologies, 551
volume 1 (long and short papers)*, pages 4171–4186. 552

Amanda Duarte, Shruti Palaskar, Lucas Ventura, Deepti 553
Ghadiyaram, Kenneth DeHaan, Florian Metze, Jordi 554
Torres, and Xavier Giro-i Nieto. 2021. How2Sign: 555
A Large-scale Multimodal Dataset for Continuous 556
American Sign Language. In *Conference on Com- 557
puter Vision and Pattern Recognition (CVPR)*. 558

Yvette Graham, Barry Haddow, and Philipp Koehn. 559
2020. [Statistical power and translationese in machine 560
translation evaluation](#). In *Proceedings of the 2020 561
Conference on Empirical Methods in Natural Lan- 562
guage Processing (EMNLP)*, pages 72–81, Online. 563
Association for Computational Linguistics. 564

Hezhen Hu, Weichao Zhao, Wengang Zhou, and 565
Houqiang Li. 2023. Signbert+: Hand-model-aware 566
self-supervised pre-training for sign language under- 567
standing. *IEEE Transactions on Pattern Analysis and 568
Machine Intelligence*, 45(9):11221–11239. 569

Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke 570
Zettlemoyer. 2018. [Adversarial example generation 571
with syntactically controlled paraphrase networks](#). In 572
*Proceedings of the 2018 Conference of the North 573
American Chapter of the Association for Computa- 574
tional Linguistics: Human Language Technologies, 575
Volume 1 (Long Papers)*, pages 1875–1885, New Or- 576
leans, Louisiana. Association for Computational Lin- 577
guistics. 578

Youngjoon Jang, Haran Raajesh, Liliane Momeni, Gül 579
Varol, and Andrew Zisserman. 2025. Lost in transla- 580
tion, found in context: Sign language translation with 581
contextual cues. *arXiv preprint arXiv:2501.09754*. 582

Vladimir I. Levenshtein. 1965. [Binary codes capable of 583
correcting deletions, insertions, and reversals](#). *Dok- 584
lady Akademii Nauk SSSR*, 163(4):845–848. 585

Zecheng Li, Wengang Zhou, Weichao Zhao, Kepeng 586
Wu, Hezhen Hu, and Houqiang Li. 2025. [Uni-sign: 587
Toward unified sign language understanding at scale](#). 588
In *The Thirteenth International Conference on Learn- 589
ing Representations*. 590

Kezhou Lin, Xiaohan Wang, Linchao Zhu, Ke Sun, 591
Yi Yang, and 1 others. 2023. Gloss-free end-to-end 592
sign language translation. In *The 61st Annual Meet- 593
ing Of The Association For Computational Linguis- 594
tics*. 595

Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris 596
McClanahan, Esha Uboweja, Michael Hays, Fan 597
Zhang, Chuo-Ling Chang, Ming Yong, Juhyun Lee, 598
and 1 others. 2019. Mediapipe: A framework for per- 599
ceiving and processing reality. In *Third workshop on 600
computer vision for AR/VR at IEEE computer vision 601
and pattern recognition (CVPR)*, volume 2019. 602

603	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Tomas Zelezny, Jakub Straka, Vaclav Javorek, Ondrej	656
604	Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	Valach, Marek Hruz, and Ivan Gruber. 2025. Exploring pose-based sign language translation: Ablation studies and attention insights. <i>arXiv preprint arXiv:2507.01532</i> .	657
605			658
606			659
607	Arindam Sengupta, Feng Jin, Renyuan Zhang, and Siyang Cao. 2020. mm-pose: Real-time human skeletal posture estimation using mmwave radars and cnns. <i>IEEE sensors journal</i> , 20(17):10032–10044.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert . <i>Preprint</i> , arXiv:1904.09675.	661
608			662
609			663
610			664
611	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	A Supplementary material	665
612			
613			
614		A.1 SLT Evaluation Details	666
615			
616		In order to preserve consistency of our experiments, we run each of our experiments with three different seeds, as discussed in Section 5. Detailed results for all of our runs and all metrics we evaluated on are reported in Table 3.	667
617			668
618	Lingfeng Shen, Lemao Liu, Haiyun Jiang, and Shuming Shi. 2022. On the evaluation metrics for paraphrase generation . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3178–3190, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		669
619			670
620			671
621			
622			
623			
624	Garrett Tanzer. 2025. FLEURS-ASL: Including American Sign Language in massively multilingual multi-task evaluation . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6167–6191, Albuquerque, New Mexico. Association for Computational Linguistics.		
625			
626			
627			
628			
629			
630			
631			
632	Garrett Tanzer and Biao Zhang. 2025. Youtube-SL-25: A large-scale, open-domain multilingual sign language parallel corpus . In <i>The Thirteenth International Conference on Learning Representations</i> .		
633			
634			
635			
636	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .		
637			
638			
639			
640			
641			
642	Dave Uthus, Garrett Tanzer, and Manfred Georg. 2023. Youtube-asl: A large-scale, open-domain american sign language-english parallel corpus . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 29029–29047. Curran Associates, Inc.		
643			
644			
645			
646			
647	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.		
648			
649			
650			
651			
652	Tomas Zelezny, Marek Hruz, Jakub Straka, and Shester Gueuwou. 2024. YouTube-ASL clip keypoint dataset . LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL).		
653			
654			
655			

Paraphrase mode	Seed	BLEU-1	BLEU-2	BLEU-3	BLEU-4	BLEURT	Rouge-L
Evaluation without paraphrases							
No paraphrases	42	28.34	16.38	10.76	7.44	0.43	0.29
No paraphrases	0	27.16	15.90	10.59	7.45	0.43	0.29
No paraphrases	1	28.69	16.50	10.83	7.48	0.42	0.29
Random	42	24.80	13.97	9.02	6.15	0.43	0.27
Random	0	26.64	14.78	9.37	6.34	0.42	0.26
Random	1	24.48	13.81	8.86	6.02	0.43	0.26
Min loss	42	26.79	15.14	9.86	6.82	0.42	0.27
Min loss	0	22.23	12.78	8.35	5.76	0.43	0.28
Min loss	1	26.87	15.41	10.03	6.89	0.43	0.27
Evaluation with paraphrases							
No paraphrases	42	33.75	19.78	12.92	8.85	0.44	0.31
No paraphrases	0	32.43	19.20	12.71	8.85	0.45	0.31
No paraphrases	1	34.06	19.91	12.99	8.88	0.44	0.31
Random	42	30.33	17.75	11.75	7.85	0.45	0.30
Random	0	33.03	19.15	12.30	8.30	0.44	0.30
Random	1	29.93	17.49	11.34	7.71	0.45	0.29
Min loss	42	32.51	18.86	12.29	8.43	0.44	0.30
Min loss	0	26.76	15.72	10.90	7.06	0.45	0.30
Min loss	1	32.42	19.09	12.44	8.51	0.44	0.30

Table 3: Performance of our models trained with different seeds and different paraphrase-based training strategies on the How2Sign test set. The checkpoints were evaluated both without and with paraphrases, where metrics are computed for all available paraphrases of each test instance and the highest-scoring paraphrase is selected. No par. denotes training with only canonical targets, Random indicates random sampling from available paraphrases during training, and Min loss selects the paraphrase yielding the minimum training loss per instance. For the Min loss configuration, a paraphrased target was selected in approximately 60% of training instances.