
SignVerse-2M: A Two-Million-Clip Pose-Native Universe of 25+ Sign Languages

Sen Fang¹, Hongbin Zhong², Yanxin Zhang³
Dimitris N. Metaxas¹

¹Rutgers University ²Georgia Institute of Technology ³NVIDIA



Dataset



Project Page



Benchmark

Abstract

Existing large-scale sign language resources typically provide supervision only at the level of raw video-text alignment and are often produced in laboratory settings. While such resources are important for semantic understanding, they are not naturally suited to open-world recognition and translation, or to modern pose-driven sign language video generation frameworks: 1. RGB-based pretrained recognition models depend heavily on fixed backgrounds or clothing conditions during recording, and are less robust in open-world settings than style-agnostic pose-processing models. 2. Recent pose-guided image/video generation models mostly use a unified keypoint representation such as DWPose as their control interface. At present, the sign language field still lacks a data resource that can directly interface with this modern pose-native paradigm while also targeting real-world open scenarios.

We present **SignVerse-2M**, a large-scale multilingual pose-native dataset for sign language pose modeling and evaluation. Built from publicly available multilingual sign language video resources, it applies DWPose in a unified preprocessing pipeline to convert raw videos into 2D pose sequences that can be used directly for modeling, resulting in a consolidated corpus of about two million clips covering more than 25 sign languages. Unlike many laboratory datasets, this resource retains the shooting conditions and speaker diversity in the real world, excluding background changes and clothing differences. Therefore, it is more suitable for researching robust sign language models in open scenarios. Toward this goal, we further provide the data construction pipeline, task definitions, and a simple SignDW Transformer baseline, demonstrating the feasibility of this resource for multilingual pose-space modeling and discussing the evaluation claims it can support as well as its current limitations.

1 Introduction

In recent years, the scale and coverage of sign language datasets have grown substantially with the expansion of online video resources and multimodal learning methods Duarte et al. [2021b], Camgöz et al. [2018], Cihan Camgöz et al. [2020], Forster et al. [2012], Joze and Koller [2018], Li et al. [2020]. However, existing large-scale sign language resources still face two main limitations: first, they are mostly designed for tasks such as sign language recognition, retrieval, or sign-to-text translation Hu et al. [2023], Jiang et al. [2021], Tarrés et al. [2023], Cihan Camgöz et al. [2020], Koller [2020], with supervision centered on weak or strong alignment between videos and text Camgöz et al. [2018], Duarte et al. [2021b], Albanie et al. [2020]; second, they often involve only a small number of signers

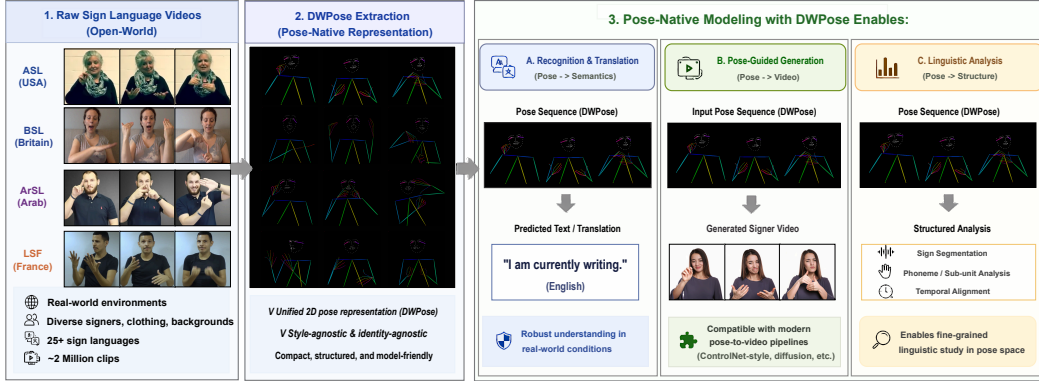


Figure 1: **Overview of SignVerse-2M.** SignVerse-2M organizes large-scale public sign language videos into a unified pose-native interface for multilingual sign language modeling. This representation is directly compatible with modern pose-driven generation pipelines and can serve as an intermediate control space for digital human or avatar generation Chen et al. [2023], Cai et al. [2023], Zwitserlood et al. [2004], Zhang and Agrawala [2023].

35 recorded with fixed clothing and backgrounds, in environments that are neither dynamic nor diverse
 36 and are overly idealized Forster et al. [2012, 2014], Von Agris and Kraiss [2010], Boháček and Hruz
 37 [2022]. These resources are valuable for understanding what is being expressed in sign language
 38 videos, but they still lack suitable data representations and evaluation interfaces for problems that
 39 demand strong robustness in real-world recognition or in generation conditioned on complex textual
 40 inputs Fang et al. [2025c,a,d]. This issue is especially pronounced in multilingual sign language
 41 settings, where raw videos contain strong nuisance factors such as background, viewpoint, signer
 42 identity, and recording conditions, while different datasets often adopt **incompatible representations,**
 43 **preprocessing pipelines, and evaluation protocols,** making model comparison and interpretation of
 44 conclusions difficult Yin et al. [2022], Rayner et al. [2016], Moryossef and Müller [2021].

45 This leads us to believe that what sign language research truly lacks today is not simply more videos,
 46 but rather a **unified, pose-driven evaluation resource for in-the-wild sign language tasks** Saunders
 47 et al. [2020b], Fang et al. [2025a], Baltatzis et al. [2024]. This need is particularly evident in sign
 48 language generation: existing models based on Mesh Human, MediaPipe, or non-standard OpenPose
 49 derivatives Cai et al. [2023], Lugaresi et al. [2019], Cao et al. [2017] are difficult to develop further
 50 with strong modern backbones released by major companies because they are not compatible with
 51 contemporary pose-to-video pipelines Ma et al. [2024], Zhang et al. [2025], Cheng et al. [2025].

*Over the past three years, a large number of pose-driven human image and video generation models have gradually converged on a unified technical interface: they first extract standardized keypoints such as DWPose Yang et al. [2023] from real-world videos, and then use these keypoints to drive ControlNet-style control modules Zhang and Agrawala [2023], Peng et al. [2025] or other pose-conditioned generation frameworks Zhang et al. [2023], Chan et al. [2019]. In other words, in the broader visual generation community, DWPose Yang et al. [2023], as a high-accuracy and high-speed 2D skeleton estimation model, has effectively become one of the **standard input representations for modern pose-driven generation.***

52
 53 In contrast, although the sign language field now has access to more and more video resources, it still
 54 lacks a standardized data interface that is directly compatible with this modern generative paradigm
 55 Stoll et al. [2020], Saunders et al. [2020b], Fang et al. [2025a]. This is particularly important for the
 56 NeurIPS 2026 Evaluations & Datasets Track: the value of a dataset should not be measured only by
 57 scale, but also by whether it supports clear, comparable, and interpretable evaluation.

58 Motivated by this perspective, we propose **SignVerse-2M.** Rather than releasing another raw video-
 59 text parallel corpus, we build on existing large-scale multilingual sign language video resources
 60 and apply DWPose for unified pose extraction and standardized processing, converting them into a
 61 large-scale pose-native sign language representation dataset covering more than 25 sign languages
 62 and about two million clips. As shown in Figure 1, our goal is not to recreate yet another video-

Dataset	Duration(h)	Vocabulary(k)	Annotation Type	Year	Domain
KETI Ko et al. [2019]	27.99	0.49	Spoken Text	2019	Emergency situations
PHOENIX-2014T Camgoz et al. [2018]	10.50	3.90	Spoken Text, Gloss	2018	Weather Forecast
CSL Daily Zhou et al. [2021]	23.27	4.60	Spoken Text, Gloss	2021	Daily life
OpenASL Shi et al. [2022]	288	33	Spoken Text	2022	Youtube (news + vlogs)
How2Sign Duarte et al. [2021a]	79.10	22.40	Spoken Text	2021	Instructional
ASLLRP Neidle et al. [2022]	80	2.1	Gloss	2022	Comprehensive
YouTube-ASL Uthus et al. [2023]	~1000	60	Spoken Text	2023	Youtube (open-domain ASL)
YouTube-SL-25 Tanzer et al. [2024]	3207	~374	Spoken Text	2024	Youtube (25+ languages)
SignVerse-2M (ours)	3207	~374	Spoken Text, Pose	2026	Pose-native open-world multilingual

Table 1: **Comparison with representative and widely used sign language datasets:** We compare SignVerse-2M with representative and widely used sign language datasets under a unified layout, covering annotation type, scale, and domain. The top five rows mainly correspond to spoken-text translation benchmarks, while ASLLRP, YouTube-ASL, and YouTube-SL-25 place our resource in broader and more open-domain settings. Duration and vocabulary are reported as corpus-level statistics; “-” denotes values that are unavailable or not directly comparable, and approximate values are marked with “~”.

63 text corpus, but to construct a data interface that can directly support sign language generation
64 research: 1. it reduces variation in visual appearance and shifts the focus back to the hand, body, and
65 upper-body motions themselves; 2. it provides a more natural input-output space for cross-lingual
66 modeling, unified training, and standardized evaluation; 3. more importantly, this representation is
67 naturally compatible with existing pose-to-pose, pose-to-video, and other pose-driven generation
68 pipelines, allowing sign language generation to move beyond starting directly from pixel space and
69 instead build on mainstream pose-conditioned modeling paradigms, benefiting from large-scale video
70 models developed by major companies. As shown in Table 1, representative existing datasets differ
71 substantially in annotation type, scale, and domain, but most are still organized around video–text
72 supervision rather than a unified pose-native interface.

73 To address these questions, this paper makes four main contributions:

- 74 • We construct **SignVerse-2M**, a large-scale, multilingual, uniformly represented pose-native
75 sign language dataset that systematically converts public video resources into a pose-
76 sequence corpus suitable for generative tasks.
- 77 • We explicitly position this resource as a **sign language interface for modern pose-driven**
78 **generative models**, emphasizing its compatibility with the DWPose-centered control
79 paradigm rather than treating pose extraction merely as a preprocessing byproduct.
- 80 • We introduce **task settings and a simple SignDW Transformer baseline** for sign lan-
81 guage generation research to validate the practical usability of this resource and provide a
82 reproducible starting point for future work.
- 83 • We discuss the **strengths** of this unified pose representation in terms of real-world robustness,
84 cross-lingual transfer, and deployment feasibility, as well as its **limitations** in fine-grained
85 hand information, non-manual features, and linguistic completeness.

86 In summary, SignVerse-2M is intended to support the shift of sign language research from idealized
87 laboratory settings toward real-world pose modeling.

88 2 Background and Positioning

89 2.1 Traditional Video-Text Sign Language Corpora

90 **The Development of Sign Language Data.** Deep learning based sign language modeling typically
91 requires large amounts of data, and in general, more data is better Camgöz et al. [2018], Cihan Camgöz
92 et al. [2020], Hu et al. [2023]. From Phoenix Forster et al. [2012, 2014], MS-ASL Joze and Koller
93 [2018], How2Sign Duarte et al. [2021b], Sign.MT, and Prompt2Sign Fang et al. [2025a] to YouTube-
94 SL-25, the overall trend has been toward either larger datasets or broader language coverage Li et al.
95 [2020], Albanie et al. [2020, 2021], Shi et al. [2022], Yin et al. [2022]. Most existing large-scale sign
96 language datasets are built around the mapping between video and text, and are primarily intended
97 for tasks such as sign language recognition, retrieval, and sign-to-text translation Cihan Camgöz et al.
98 [2020], Tarrés et al. [2023], Jiang et al. [2021], Koller [2020]. For these tasks, raw video preserves

99 the richest visual information, while text provides clear semantic supervision, making video-text
100 parallel corpora the natural data format. However, when the research goal shifts to translation in
101 open environments or to modern sign language video generation Saunders et al. [2020b], Fang et al.
102 [2025a], Baltatzis et al. [2024], Yin et al. [2024], the problem changes. The focus is no longer only
103 on whether semantics are aligned, but also on whether long motion sequences are coherent, whether
104 hand and upper-body movements are natural, whether representations are consistent across samples,
105 and whether different models can be fairly compared under a unified interface Stoll et al. [2020],
106 Saunders et al. [2020b], Walsh et al. [2025]. In this setting, directly using raw video often mixes
107 irrelevant factors such as background, signer appearance, camera viewpoint, and recording conditions
108 into both modeling and evaluation, thereby weakening the analysis of motion itself Boháček and
109 Hrúz [2022], Fang et al. [2025d].

110 2.2 From Laboratory Distributions to Real-World Distributions

111 **Why Choose a Pose-Native Representation?** Pose representation provides a more suitable inter-
112 mediate space for this problem. By uniformly mapping raw videos into keypoint sequences Yang
113 et al. [2023], Cao et al. [2017], Lugaresi et al. [2019], Cai et al. [2023], models can focus more
114 directly on the temporal structure and body geometry of sign language motion without having to
115 simultaneously handle complex appearance variation Saunders et al. [2020b, 2021], Stoll et al. [2020].
116 This is especially important in multilingual settings, where the visual differences across languages,
117 video sources, and signers are often much larger than the domain gaps commonly seen in text corpora
118 Yin et al. [2022], Fang et al. [2025a]. More importantly, over the past three years, DWPose Yang et al.
119 [2023] has become one of the standardized control interfaces for a large number of pose-driven image
120 and video generation models Ma et al. [2024], Zhang et al. [2025], Cheng et al. [2025]. Whether in
121 ControlNet-style conditional control Zhang and Agrawala [2023], Zhang et al. [2023], Peng et al.
122 [2025] or in other generation frameworks that use human keypoints as intermediate conditions Chan
123 et al. [2019], Saunders et al. [2020a], a stable, general, and large-scale extractable pose representation
124 is essential. A concrete example is that RGB-based recognition models are difficult to apply directly
125 in real-world environments: once the background, person, or clothing changes, the input video
126 features can differ significantly Koller [2020], Boháček and Hrúz [2022]. Pose-based inputs or
127 generation, in contrast, are more robust Fang et al. [2025c,b].

128 2.3 Positioning of SignVerse-2M

129 **Relation to Existing Multilingual Sign Language Resources.** SignVerse-2M inherits from exist-
130 ing public multilingual sign language video corpora Duarte et al. [2021b], Albanie et al. [2020, 2021],
131 Shi et al. [2022], Yin et al. [2022], but its research objective is different. Prior work primarily asks
132 how to collect, align, and organize cross-lingual sign language video-text data from the web in order
133 to support translation tasks Cihan Camgöz et al. [2020], Camgöz et al. [2018]. SignVerse-2M instead
134 focuses on a different question: when the target setting shifts toward the real world, can these public
135 video resources be transformed into a unified pose space that serves as reusable, extensible, and
136 comparable infrastructure for generative research Yang et al. [2023], Cheng et al. [2025], Zhang and
137 Agrawala [2023]? We position SignVerse-2M as a pose-native resource not because pose extraction
138 itself is especially complex, but because it enables sign language data to connect naturally with
139 mainstream generative paradigms Ma et al. [2024], Zhang et al. [2025], Peebles and Xie [2023],
140 making the definition of sign language motion representation and its associated evaluation protocol
141 part of the scientific object itself rather than merely a preprocessing step Fang et al. [2025a]. We
142 therefore hope that reviewers and readers will understand this work as a reconstruction of the *data*
143 *representation, task interface, and evaluation object*, rather than as a simple expansion of existing
144 video-text corpora.

145 **Questions Addressed in This Paper.** Based on the above considerations, this paper does not aim
146 to prove that SignVerse-2M is superior to existing resources for every sign language task. Instead, it
147 focuses on three more specific questions. First, can large-scale public videos be stably transformed
148 into a unified pose corpus suitable for sign language generation and directly compatible with modern
149 pose-driven generation frameworks Yang et al. [2023], Cheng et al. [2025], Ma et al. [2024]? Second,
150 is this pose-space representation sufficient to support the training and comparison of multilingual
151 generative models under real-world distributions Yin et al. [2022], Fang et al. [2025a,c]? Third, if
152 used as an evaluation resource, what conclusions can it support, and where are its boundaries and

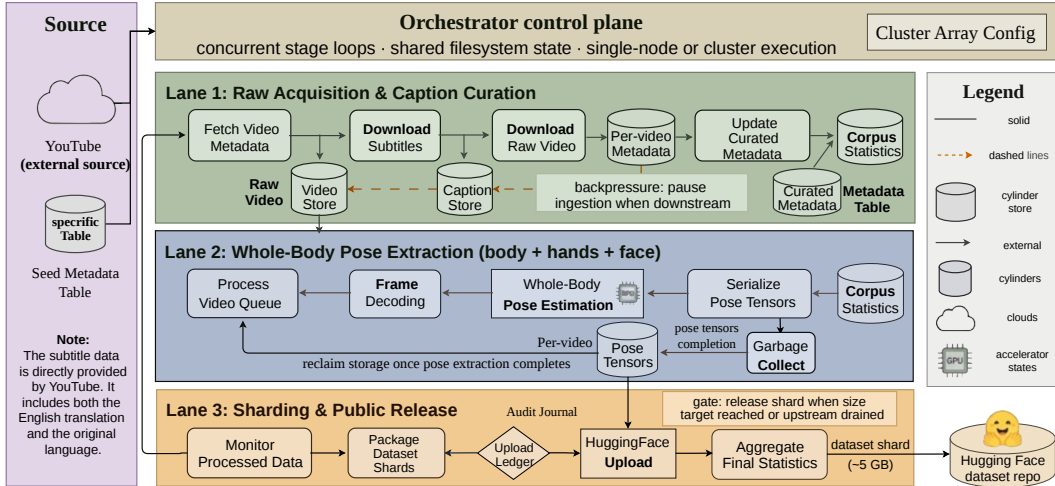


Figure 2: **Overview of the SignVerse-2M data processing pipeline.** The pipeline is organized into three main lanes: raw acquisition and caption curation, whole-body pose extraction, and sharding for public release. Starting from a manifest indexed by ‘video_id’ and ‘sign_language’, the system retrieves metadata, subtitles, and raw videos, converts each video into DWPose-based body, hand, and face keypoint sequences, and then packages the processed outputs into dataset shards for publication. The orchestration layer maintains status records, supports staged re-execution and failure recovery, and enables the corpus to be built incrementally rather than as a one-shot preprocessing export.

153 limitations Walsh et al. [2025], Koller [2020]? The remainder of the paper is organized around these
 154 three questions.

155 3 Data Infrastructure and Pipeline

156 3.1 Data Source and Curation Pipeline

157 **Data Source.** SignVerse-2M starts from publicly available multilingual sign language video re-
 158 sources and builds a large-scale corpus for pose-based modeling on top of them. We maintain
 159 a manifest centered on ‘video_id’ and ‘sign_language’, which drives the subsequent processing
 160 pipeline. The raw metadata table in the current repository contains 39,196 video entries and roughly
 161 2M segments; each entry is fed into a traceable, resumable, and re-runnable processing pipeline with
 162 associated status records. Accordingly, the focus of SignVerse-2M is to build an infrastructure for
 163 automatic sign language video processing that turns open-world videos into a unified, trainable, and
 164 publishable pose-native corpus.

165 **Pipeline Overview.** Figure 2 summarizes the overall construction path of SignVerse-2M. The
 166 pipeline begins with manifest-driven sample management, proceeds through metadata and subtitle
 167 acquisition, and then moves into pose extraction, packaging, and publication. This layout is useful
 168 because it makes clear that the corpus is not defined by a single preprocessing script, but by a sequence
 169 of recoverable and independently executable stages.

170 **Video Acquisition.** Before pose extraction begins, we first complete video-level data acquisition.
 171 For each video ID in the manifest, the system retrieves the original platform metadata, downloads
 172 the raw video, and fetches the available subtitles in parallel. This stage continuously writes back
 173 fields such as the title, duration, subtitle languages, raw video path, processing time, and error status,
 174 allowing the metadata table to serve as both a sample inventory and a runtime record. As a result, the
 175 corpus construction process can proceed incrementally by entry and can resume from intermediate
 176 states after interruption.

177 **Subtitle Structuring and Language Signals.** Once subtitles are downloaded locally, they are
 178 converted into structured text objects instead of preserving the platform-exported format verbatim.


```

poses_schema_part1.txt
1  poses.npz
2  |-- video_id
3  |-- fps
4  |-- total_frames
5  |-- frame_indices
6  |-- frame_width, frame_height
7  |-- frame_payloads[t]
8     |-- num_persons
9     |-- frame_width, frame_height
10    |-- person_k

poses_schema_part2.txt
1  person_k
2  |-- body_keypoints[18, 2]
3  |-- body_scores[18]
4  |-- face_keypoints[68, 2]
5  |-- face_scores[68]
6  |-- left_hand_keypoints[21, 2]
7  |-- left_hand_scores[21]
8  |-- right_hand_keypoints[21, 2]
9  |-- right_hand_scores[21]

```

Figure 4: **Schema of the released ‘poses.npz’ payload.** The stored format is person-centric: each frame payload records the number of detected signers and organizes body, face, left-hand, and right-hand keypoints together with their confidence scores under each ‘person_k’. This structure is the native representation released by SignVerse-2M and is the basis from which visualization scripts derive OpenPose-style aggregated renderings.

218 separated into distinct stages and can be executed on selected subsets of samples, which leaves room
 219 for future expansion and maintenance.

220 **Cross-Lingual Scale.** Figure 3 gives a direct view of the language distribution inherited from
 221 YouTube-SL-25. The logarithmic y-axis makes the long-tail structure visible: a small number of
 222 languages contribute a large fraction of the total hours, while many other languages remain present at
 223 much smaller scales. For SignVerse-2M, this distribution is not a side detail but a core property of
 224 the corpus, because it shapes multilingual training, transfer behavior, and the practical difficulty of
 225 benchmarking across languages.

226 4.2 DWPose and Visualization

227 **DWPose Representation and Visualization.** SignVerse-2M uses DWPose as the backend for pose
 228 extraction. For each frame, the system extracts body, hand, and facial keypoints together with their
 229 confidence scores. In the aggregated format, ‘poses.npz’ stores not only video-level metadata but also
 230 per-frame payloads; each payload records the number of people, frame dimensions, and person-wise
 231 body, face, left-hand, and right-hand keypoints and scores. This representation allows downstream
 232 methods to consume temporal inputs directly at the keypoint level without re-decoding the original
 233 videos. The dataset repository also provides visualization scripts that reconstruct ‘poses.npz’ into
 234 inspectable skeletal renderings. Specifically, the scripts can read either aggregated or per-frame NPZ
 235 outputs and generate multiple visualization styles, including control-style skeletons, OpenPose-style
 236 renderings, and previews overlaid on the original video. The underlying drawing functions render
 237 body, hand, and face keypoints separately while using confidence scores to modulate the display.
 238 As shown in Figure 4, the released ‘poses.npz’ format stores pose data in a person-centric structure
 239 rather than as a single flat OpenPose-style array. These visualization interfaces are not intended as
 240 part of the training input itself; instead, they serve as inspection and debugging tools for quickly
 241 identifying missing keypoints, multi-person interference, local jitter, and abnormal frames.

242 **Processing Skeleton.** From an implementation perspective, the more informative summary is the
 243 structure of the stored pose payload rather than another pipeline sketch. The released ‘poses.npz’
 244 format can be summarized by the following schema-like pseudocode:

245 This organization makes explicit that SignVerse-2M stores pose data in a person-centric DWPose-
 246 style payload, with body, face, and both hands grouped under each detected signer at each frame.
 247 The accompanying visualization scripts can reorganize the same information into an OpenPose-style
 248 aggregated representation for rendering, but the released corpus itself keeps the per-person structure
 249 because it is easier to inspect, serialize, and reuse in downstream modeling.

Model	Language	DEV SET					TEST SET				
		BLEU-4	BLEU-3	BLEU-2	BLEU-1	ROUGE	BLEU-4	BLEU-3	BLEU-2	BLEU-1	ROUGE
SignDW Transformer (40M)	ASL	11.27	17.11	25.82	34.97	35.05	9.65	15.51	24.40	33.85	33.71
SignDW Transformer (1.2B)	ASL	11.72	17.56	26.12	35.82	35.89	10.10	15.96	24.70	34.30	34.26
SignDW Transformer (40M)	BSL	9.58	15.32	24.05	33.80	33.42	8.35	14.04	22.72	32.58	32.08
SignDW Transformer (1.2B)	BSL	10.06	16.06	25.07	35.02	34.77	8.80	14.49	23.12	33.23	33.03
SignDW Transformer (40M)	GSL	9.82	15.68	24.48	34.23	33.98	9.69	15.48	24.29	33.82	33.58
SignDW Transformer (1.2B)	GSL	10.70	16.80	25.94	35.73	35.55	10.14	16.19	25.28	34.97	35.04
SignDW Transformer (40M)	DSGS	10.16	16.07	25.02	34.73	34.25	9.12	14.96	23.77	33.10	33.09
SignDW Transformer (1.2B)	DSGS	10.70	16.70	25.61	35.49	35.56	9.57	15.41	24.07	33.73	33.87

Table 2: **Text-to-pose baseline results on four sign languages.** We report results for **SignDW Transformer (40M)** and **SignDW Transformer (1.2B)** on ASL, BSL, GSL, and DSGS.

250 5 Evaluation & Baseline

251 Given the constraints of space and computational resources, we focus in this paper on a sign language
 252 generation baseline, as this setting more directly highlights the strengths of our dataset.

253 5.1 Setup

254 **Task.** Sign language generation aims to synthesize sign language motion from complex textual inputs.
 255 In this paper, our quantitative evaluation is carried out in pose space: the generated DWPose sequence
 256 is translated back into spoken-language text, and the recovered sentence is then compared with the
 257 original input. This back-translation protocol provides a broad proxy for semantic fidelity while
 258 avoiding additional variation introduced by downstream video rendering Saunders et al. [2020b].

259 We evaluate the pipeline at two levels. First, we measure the performance of the text-to-DWPose
 260 model itself using automatic back-translation metrics. Second, we feed the generated pose videos into
 261 modern pose-driven rendering models to synthesize final videos and inspect the overall generation
 262 quality qualitatively. The rendered videos are therefore used as a compatibility and visualization
 263 check rather than as part of the automatic evaluation protocol.

264 For pose-space back-translation, we use SignX Fang et al. [2026], a recent sign language translation
 265 model trained on pose data in the corresponding language. Depending on the dataset, its BLEU-4
 266 performance typically falls in the range of 25–28, making it a strong translation model for this
 267 evaluation setting.

268 **Metrics.** (i) BLEU- n score measures the similarity between the generated translation and the
 269 reference text based on n -gram overlap. Higher scores indicate that the prediction is closer to the
 270 reference sentence. Larger values of n impose stricter requirements on local fluency and phrase-level
 271 consistency Papineni et al. [2002].

272 (ii) ROUGE score Lin [2004] is similar in spirit to BLEU, but places greater emphasis on coverage
 273 and overlap with the reference text. A higher ROUGE score indicates that the generated text is more
 274 consistent with the reference and is therefore more complete and accurate.

275 5.2 Results.

276 Table 2 summarizes the text-to-pose results on four representative sign languages: ASL, BSL, GSL,
 277 and DSGS. For all runs, we train the baseline for 200K steps with batch size 32 and an initial learning
 278 rate of 0.001 under the same multilingual DWPose interface, so that differences across rows mainly
 279 reflect model scale and language-specific data conditions rather than changes in the optimization
 280 setting. Overall, the proposed baseline achieves stable performance across all four settings, which
 281 suggests that the unified DWPose representation is sufficient to support multilingual text-to-pose
 282 modeling under a shared training and evaluation interface. Across languages, the larger **SignDW**
 283 **Transformer (1.2B)** model is consistently stronger than the **SignDW Transformer (40M)** model,
 284 with the clearest gains appearing in BLEU-4 and BLEU-3, indicating better preservation of longer
 285 local motion patterns after back-translation. At the same time, the gap between languages remains
 286 non-negligible: ASL and GSL are relatively stronger, while BSL and DSGS remain more challenging,
 287 especially on the test split. We view this pattern as consistent with the multilingual and open-world
 288 nature of SignVerse-2M. The benchmark is therefore not only a proof of feasibility, but also a first

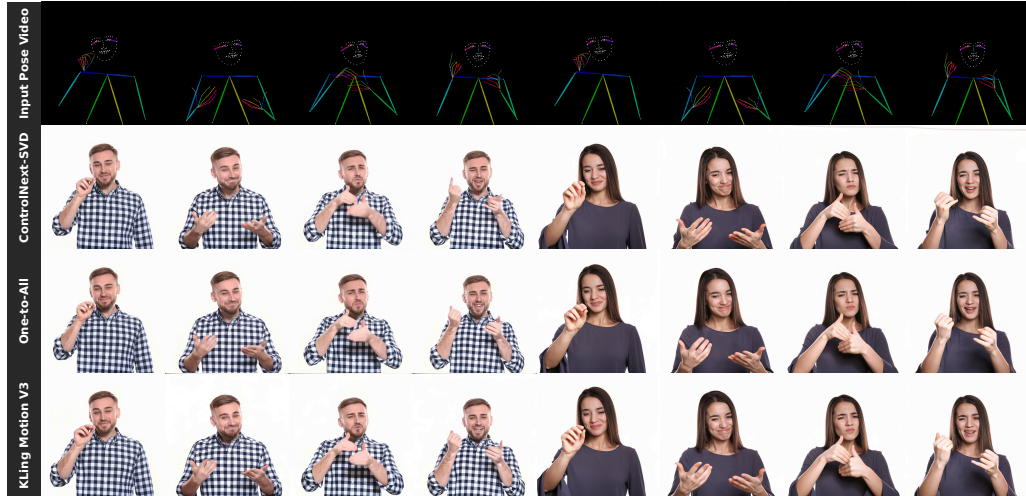


Figure 5: **Qualitative comparison of pose-conditioned sign video rendering.** The first row shows the input DWPose video. To remain compatible with the model-specific preprocessing pipelines of different renderers, we prepend a body-shape mask before feeding the pose video into each method. The next three rows compare **ControlNext-SVD** Peng et al. [2025], **One-to-All** Shi et al. [2025], and **KLing Motion V3** Team et al. [2025] under the same pose sequence. For each rendering row, the left half shows the male signer and the right half shows the female signer.

289 indication that model scale and language-specific data conditions both matter under a pose-native
 290 sign generation setting.

291 Figure 5 provides a complementary qualitative view of this result. Beyond back-translation scores,
 292 the examples show that the DWPose-based representation used in SignVerse-2M can be interfaced
 293 directly with several modern pose-conditioned rendering systems, including both commercial and
 294 open-source pipelines, while still preserving the signer’s coarse motion structure and temporal
 295 progression. This point is important for our dataset motivation: earlier sign-language pipelines
 296 often relied on representations such as MediaPipe pose, OpenPose-style subsets, or SMPL-derived
 297 intermediate formats, which are useful for analysis but are not naturally aligned with the control
 298 interfaces adopted by recent image-to-video and talking-human generation models. In contrast,
 299 the representation used here makes the downstream rendering step substantially easier to connect,
 300 compare, and reuse, which is precisely the kind of systems-level compatibility that a pose-native
 301 benchmark should provide.

302 6 Conclusion

303 We present SignVerse-2M, a large-scale multilingual pose-native dataset for sign language research.
 304 Unlike existing sign language resources that are primarily derived from laboratory settings, the core
 305 contribution of SignVerse-2M is neither simply to scale up the amount of video data nor to introduce
 306 a new pose estimation algorithm. Instead, it systematically converts publicly available multilingual
 307 sign language videos from open environments into a unified DWPose representation, and builds
 308 around this representation the data interface and evaluation foundation needed for sign language tasks.
 309 Through this resource, we aim to shift the focus of sign language generation research away from
 310 heterogeneous RGB video preprocessing and incomparable experimental settings toward motion
 311 representation, cross-lingual transfer, and evaluation itself.

312 Overall, we hope that SignVerse-2M will provide not merely a collection of samples, but a new
 313 research interface. This resource helps move the community’s focus from performance validation in
 314 idealized laboratory settings toward robust modeling and deployment evaluation under real-world
 315 conditions, thereby offering a more unified, transparent, and comparable foundation for future sign
 316 language research.

317 **References**

- 318 Samuel Albanie, Gül Varol, Liliane Momeni, Triantafyllos Afouras, Joon Son Chung, Neil Fox, and Andrew
319 Zisserman. BSL-1K: Scaling up Co-articulated Sign Language Recognition using Mouthing Cues. In
320 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- 321 Samuel Albanie, Gül Varol, Liliane Momeni, Hannah Bull, Triantafyllos Afouras, Himel Chowdhury, Neil Fox,
322 Bencie Woll, Rob Cooper, Andrew McParland, and Andrew Zisserman. BOBSL: BBC-Oxford British Sign
323 Language Dataset. 2021.
- 324 Vasileios Baltatzis, Rolandos Alexandros Potamias, Evangelos Ververas, Guanxiong Sun, Jiankang Deng, and
325 Stefanos Zafeiriou. Neural sign actors: A diffusion model for 3d sign language production from text, 2024.
326 URL <https://arxiv.org/abs/2312.02702>.
- 327 Matyáš Boháček and Marek Hruží. Sign pose-based transformer for word-level sign language recognition. In
328 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops*,
329 pages 182–191, January 2022.
- 330 Zhongang Cai, Wanqi Yin, Ailing Zeng, Chen Wei, Qingping Sun, Yanjun Wang, Hui En Pang, Haiyi Mei,
331 Mingyuan Zhang, Lei Zhang, Chen Change Loy, Lei Yang, and Ziwei Liu. Smler-x: Scaling up expressive
332 human pose and shape estimation, 2023.
- 333 Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. Neural sign language
334 translation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7784–7793,
335 2018. doi: 10.1109/CVPR.2018.00812.
- 336 Necati Cihan Camgöz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. Neural Sign
337 Language Translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
338 *(CVPR)*, 2018.
- 339 Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. OpenPose: Realtime Multi-Person 2D
340 Pose Estimation using Part Affinity Fields. In *Proceedings of the IEEE Conference on Computer Vision and*
341 *Pattern Recognition (CVPR)*, 2017.
- 342 Caroline Chan, Shiry Ginossar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now, 2019. URL
343 <https://arxiv.org/abs/1808.07371>.
- 344 Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands
345 via motion diffusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
346 *Pattern Recognition*, pages 18000–18010, 2023.
- 347 Gang Cheng, Xin Gao, Li Hu, Siqi Hu, Mingyang Huang, Chaonan Ji, Ju Li, Dechao Meng, Jinwei Qi, Penchong
348 Qiao, et al. Wan-animate: Unified character animation and replacement with holistic replication. *arXiv*
349 *preprint arXiv:2509.14055*, 2025.
- 350 Necati Cihan Camgöz, Oscar Koller, Simon Hadfield, and Richard Bowden. Sign Language Transformers:
351 Joint End-to-end Sign Language Recognition and Translation. In *Proceedings of the IEEE Conference on*
352 *Computer Vision and Pattern Recognition (CVPR)*, 2020.
- 353 Amanda Duarte, Shruti Palaskar, Lucas Ventura, Deepti Ghadiyaram, Kenneth DeHaan, Florian Metze, Jordi
354 Torres, and Xavier Giro-i Nieto. How2Sign: A Large-scale Multimodal Dataset for Continuous American
355 Sign Language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
356 *(CVPR)*, 2021a.
- 357 Amanda Duarte, Shruti Palaskar, Lucas Ventura, Deepti Ghadiyaram, Kenneth DeHaan, Florian Metze, Jordi
358 Torres, and Xavier Giro-i Nieto. How2Sign: A Large-Scale Multimodal Dataset for Continuous American
359 Sign Language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
360 *(CVPR)*, 2021b.
- 361 Sen Fang, Chen Chen, Lei Wang, Ce Zheng, Chunyu Sui, and Yapeng Tian. Signllm: Sign language production
362 large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*
363 *(ICCV) Workshops*, pages 6622–6634, October 2025a.
- 364 Sen Fang, Yalin Feng, Hongbin Zhong, Yanxin Zhang, and Dimitris N. Metaxas. Stable signer: Hierarchical
365 sign language generative model, 2025b. URL <https://arxiv.org/abs/2512.04048>.
- 366 Sen Fang, Chunyu Sui, Hongwei Yi, Carol Neidle, and Dimitris N. Metaxas. Signx: The foundation model for
367 sign recognition, 2025c. URL <https://arxiv.org/abs/2504.16315>.

- 368 Sen Fang, Chunyu Sui, Yanghao Zhou, Xuedong Zhang, Hongbin Zhong, Yapeng Tian, and Chen Chen. Signdiff:
369 Diffusion model for american sign language production, 2025d. URL <https://arxiv.org/abs/2308.16082>.
370
- 371 Sen Fang, Yalin Feng, Chunyu Sui, Hongbin Zhong, Hongwei Yi, and Dimitris N. Metaxas. Signx: Continuous
372 sign recognition in compact pose-rich latent space, 2026. URL <https://arxiv.org/abs/2504.16315>.
- 373 Jens Forster, Christoph Schmidt, Thomas Hoyoux, Oscar Koller, Uwe Zelle, Justus H Piater, and Hermann Ney.
374 RWTH-PHOENIX-Weather: A Large Vocabulary Sign Language Recognition and Translation Corpus. In
375 *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, 2012.
- 376 Jens Forster, Christoph Schmidt, Oscar Koller, Martin Bellgardt, and Hermann Ney. Extensions of the sign
377 language recognition and translation corpus rwth-phoenix-weather. In *Proceedings of the Ninth International
378 Conference on Language Resources and Evaluation (LREC'14)*, pages 1911–1916, 2014.
- 379 Lianyu Hu, Liqing Gao, Zekang Liu, and Wei Feng. Continuous sign language recognition with correlation
380 network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- 381 Songyao Jiang, Bin Sun, Lichen Wang, Yue Bai, Kunpeng Li, and Yun Fu. Skeleton aware multi-modal
382 sign language recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
383 Recognition (CVPR) Workshops*, 2021.
- 384 Hamid Reza Vaezi Joze and Oscar Koller. Ms-asl: A large-scale data set and benchmark for understanding
385 american sign language. *arXiv preprint arXiv:1812.01053*, 2018.
- 386 Sang-Ki Ko, Chang Jo Kim, Hyedong Jung, and Choongsang Cho. Neural sign language translation based on
387 human keypoint estimation. *Applied sciences*, 9(13):2683, 2019.
- 388 Oscar Koller. Quantitative Survey of the State of the Art in Sign Language Recognition. *arXiv preprint
389 arXiv:2008.09918*, 2020.
- 390 Dongxu Li, Cristian Rodriguez, Xin Yu, and Hongdong Li. Word-level deep sign language recognition from
391 video: A new large-scale dataset and methods comparison. In *Proceedings of the IEEE/CVF winter conference
392 on applications of computer vision*, pages 1459–1469, 2020.
- 393 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches
394 Out*, pages 74–81. Association for Computational Linguistics, July 2004.
- 395 Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays, Fan Zhang,
396 Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework for building perception
397 pipelines. *arXiv preprint arXiv:1906.08172*, 2019.
- 398 Yue Ma, Yingqing He, Xiaodong Cun, Xintao Wang, Siran Chen, Ying Shan, Xiu Li, and Qifeng Chen. Follow
399 your pose: Pose-guided text-to-video generation using pose-free videos, 2024. URL <https://arxiv.org/abs/2304.01186>.
400
- 401 Amit Moryossef and Mathias Müller. Sign language datasets. [https://github.com/
402 sign-language-processing/datasets](https://github.com/sign-language-processing/datasets), 2021.
- 403 Carol Neidle, Augustine Opoku, and Dimitris Metaxas. ASL Video Corpora & Sign Bank: Resources Available
404 through the American Sign Language Linguistic Research Project (ASLLRP), 2022.
- 405 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation
406 of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational
407 Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational
408 Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- 409 William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL [https://arxiv.
410 org/abs/2212.09748](https://arxiv.org/abs/2212.09748).
- 411 Bohao Peng, Jian Wang, Yuechen Zhang, Wenbo Li, Ming-Chang Yang, and Jiaya Jia. Controlnext: Powerful
412 and efficient control for image and video generation, 2025. URL <https://arxiv.org/abs/2408.06070>.
- 413 Manny Rayner, Pierrette Bouillon, Sarah Ebling, Johanna Gerlach, Irene Strasly, and Nikos Tsourakis. An
414 open web platform for rule-based speech-to-sign translation. In *54th Annual Meeting of the Association for
415 Computational Linguistics*, volume 2, pages 162–168, 2016.
- 416 Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. Everybody sign now: Translating spoken language
417 to photo realistic sign language video, 2020a. URL <https://arxiv.org/abs/2011.09846>.

- 418 Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Progressive Transformers for End-to-End Sign
419 Language Production. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020b.
- 420 Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Continuous 3D Multi-Channel Sign Language
421 Production via Progressive Transformers and Mixture Density Networks. *International Journal of Computer
422 Vision (IJCV)*, 2021.
- 423 Bowen Shi, Diane Brentari, Greg Shakhnarovich, and Karen Livescu. Open-domain sign language translation
424 learned from online video. *Proceedings of the 2022 Conference on Empirical Methods in Natural Language
425 Processing*, 2022.
- 426 Shijun Shi, Jing Xu, Zhihang Li, Chunli Peng, Xiaoda Yang, Lijing Lu, Kai Hu, and Jiangning Zhang. One-to-all
427 animation: Alignment-free character animation and image pose transfer, 2025. URL [https://arxiv.org/
428 abs/2511.22940](https://arxiv.org/abs/2511.22940).
- 429 Stephanie Stoll, Necati Cihan Camgöz, Simon Hadfield, and Richard Bowden. Text2Sign: Towards Sign
430 Language Production using Neural Machine Translation and Generative Adversarial Networks. *International
431 Journal of Computer Vision (IJCV)*, 2020.
- 432 Garrett Tanzer, Biao Zhang, David Uthus, and Manfred Georg. YouTube-SL-25: A large-scale, open-domain
433 multilingual sign language parallel corpus, 2024.
- 434 Laia Tarrés, Gerard I. Gállego, Amanda Duarte, Jordi Torres, and Xavier Giró i Nieto. Sign language translation
435 from instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
436 Recognition (CVPR) Workshops*, 2023.
- 437 Kling Team, Jialu Chen, Yuanzheng Ci, Xiangyu Du, Zipeng Feng, Kun Gai, Sainan Guo, Feng Han, Jingbin
438 He, Kang He, Xiao Hu, Xiaohua Hu, Boyuan Jiang, Fangyuan Kong, Hang Li, Jie Li, Qingyu Li, Shen
439 Li, Xiaohan Li, Yan Li, Jiajun Liang, Borui Liao, Yiqiao Liao, Weihong Lin, Quande Liu, Xiaokun Liu,
440 Yilun Liu, Yuliang Liu, Shun Lu, Hangyu Mao, Yunyao Mao, Haodong Ouyang, Wenyu Qin, Wanqi Shi,
441 Xiaoyu Shi, Lianghao Su, Haozhi Sun, Peiqin Sun, Pengfei Wan, Chao Wang, Chenyu Wang, Meng Wang,
442 Qiulin Wang, Runqi Wang, Xintao Wang, Xuebo Wang, Zekun Wang, Min Wei, Tiancheng Wen, Guohao
443 Wu, Xiaoshi Wu, Zhenhua Wu, Da Xie, Yingtong Xiong, Yulong Xu, Sile Yang, Zikang Yang, Weicai
444 Ye, Ziyang Yuan, Shenglong Zhang, Shuaiyu Zhang, Yuanxing Zhang, Yufan Zhang, Wenzheng Zhao,
445 Ruijiang Zhou, Yan Zhou, Guosheng Zhu, and Yongjie Zhu. Kling-omni technical report, 2025. URL
446 <https://arxiv.org/abs/2512.16776>.
- 447 David Uthus, Garrett Tanzer, and Manfred Georg. YouTube-ASL: A large-scale, open-domain american sign
448 language-english parallel corpus, 2023.
- 449 U. Von Agris and K.-F. Kraiss. Signum database: Video corpus for signer-independent continuous sign language
450 recognition. In *Workshop on Representation and Processing of Sign Languages*, pages 243–246, 2010.
- 451 Harry Walsh, Ed Fish, Ozge Mercanoglu Sincan, Mohamed Ilyes Lakkhal, Richard Bowden, Neil Fox, Bencie
452 Woll, Kepeng Wu, Zecheng Li, Weichao Zhao, Haodong Wang, Wengang Zhou, Houqiang Li, Shengeng
453 Tang, Jiayi He, Xu Wang, Ruobei Zhang, Yaxiong Wang, Lechao Cheng, Meryem Tasyurek, Tugce Kiziltepe,
454 and Hacer Yalim Keles. Slrtp2025 sign language production challenge: Methodology, results, and future
455 work, 2025. URL <https://arxiv.org/abs/2508.06951>.
- 456 Zhendong Yang, Ailing Zeng, Chun Yuan, and Yu Li. Effective whole-body pose estimation with two-stages
457 distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4210–4220,
458 2023.
- 459 Aoxiong Yin, Zhou Zhao, Weike Jin, Meng Zhang, Xingshan Zeng, and Xiaofei He. Mslt: Towards multilingual
460 sign language translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
461 Recognition*, pages 5109–5119, 2022.
- 462 Aoxiong Yin, Haoyuan Li, Kai Shen, Siliang Tang, and Yueting Zhuang. T2s-gpt: Dynamic vector quantization
463 for autoregressive sign language production from text, 2024. URL <https://arxiv.org/abs/2406.07119>.
- 464 Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023.
- 465 Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Controlvideo:
466 Training-free controllable text-to-video generation, 2023. URL <https://arxiv.org/abs/2305.13077>.
- 467 Yuang Zhang, Jiayi Gu, Li-Wen Wang, Han Wang, Junqi Cheng, Yuefeng Zhu, and Fangyuan Zou. Mimicmotion:
468 High-quality human motion video generation with confidence-aware pose guidance, 2025. URL <https://arxiv.org/abs/2406.19680>.
- 469

- 470 Hao Zhou, Wen gang Zhou, Weizhen Qi, Junfu Pu, and Houqiang Li. Improving sign language translation with
471 monolingual data by sign back-translation. *Proceedings of the IEEE/CVF Conference on Computer Vision
472 and Pattern Recognition (CVPR)*, pages 1316–1325, 2021.
- 473 Inge Zwitterlood, Margriet Verlinden, Johan Ros, and Sanny Van Der Schoot. Synthetic Signing for the Deaf:
474 Esign. In *Proceedings of the Conference and Workshop on Assistive Technologies for Vision and Hearing
475 Impairment (CVHI)*, 2004.

476 **NeurIPS Paper Checklist**

477 **1. Claims**

478 Question: Do the main claims made in the abstract and introduction accurately reflect the
479 paper’s contributions and scope?

480 Answer: **[Yes]**

481 Justification: The abstract and introduction explicitly state the paper’s scope as a multilingual
482 pose-native dataset, its task interface, and a simple baseline; see the abstract and Sections 1
483 and 2.

484 Guidelines:

- 485 • The answer **[N/A]** means that the abstract and introduction do not include the claims
486 made in the paper.
- 487 • The abstract and/or introduction should clearly state the claims made, including the
488 contributions made in the paper and important assumptions and limitations. A **[No]** or
489 **[N/A]** answer to this question will not be perceived well by the reviewers.
- 490 • The claims made should match theoretical and experimental results, and reflect how
491 much the results can be expected to generalize to other settings.
- 492 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
493 are not attained by the paper.

494 **2. Limitations**

495 Question: Does the paper discuss the limitations of the work performed by the authors?

496 Answer: **[Yes]**

497 Justification: The paper discusses limitations of the pose-native representation, including
498 fine-grained hand information, non-manual features, and linguistic completeness; see the
499 abstract, Section 1, and the concluding discussion in Section 5.

500 Guidelines:

- 501 • The answer **[N/A]** means that the paper has no limitation while the answer **[No]** means
502 that the paper has limitations, but those are not discussed in the paper.
- 503 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 504 • The paper should point out any strong assumptions and how robust the results are to
505 violations of these assumptions (e.g., independence assumptions, noiseless settings,
506 model well-specification, asymptotic approximations only holding locally). The authors
507 should reflect on how these assumptions might be violated in practice and what the
508 implications would be.
- 509 • The authors should reflect on the scope of the claims made, e.g., if the approach was
510 only tested on a few datasets or with a few runs. In general, empirical results often
511 depend on implicit assumptions, which should be articulated.
- 512 • The authors should reflect on the factors that influence the performance of the approach.
513 For example, a facial recognition algorithm may perform poorly when image resolution
514 is low or images are taken in low lighting. Or a speech-to-text system might not be
515 used reliably to provide closed captions for online lectures because it fails to handle
516 technical jargon.
- 517 • The authors should discuss the computational efficiency of the proposed algorithms
518 and how they scale with dataset size.
- 519 • If applicable, the authors should discuss possible limitations of their approach to
520 address problems of privacy and fairness.
- 521 • While the authors might fear that complete honesty about limitations might be used by
522 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
523 limitations that aren’t acknowledged in the paper. The authors should use their best
524 judgment and recognize that individual actions in favor of transparency play an impor-
525 tant role in developing norms that preserve the integrity of the community. Reviewers
526 will be specifically instructed to not penalize honesty concerning limitations.

527 **3. Theory assumptions and proofs**

528 Question: For each theoretical result, does the paper provide the full set of assumptions and
529 a complete (and correct) proof?

530 Answer: [N/A]

531 Justification: The paper does not present formal theoretical results or proofs.

532 Guidelines:

- 533 • The answer [N/A] means that the paper does not include theoretical results.
- 534 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
535 referenced.
- 536 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 537 • The proofs can either appear in the main paper or the supplemental material, but if
538 they appear in the supplemental material, the authors are encouraged to provide a short
539 proof sketch to provide intuition.
- 540 • Inversely, any informal proof provided in the core of the paper should be complemented
541 by formal proofs provided in appendix or supplemental material.
- 542 • Theorems and Lemmas that the proof relies upon should be properly referenced.

543 4. Experimental result reproducibility

544 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
545 perimental results of the paper to the extent that it affects the main claims and/or conclusions
546 of the paper (regardless of whether the code and data are provided or not)?

547 Answer: [Yes]

548 Justification: The paper describes the data construction pipeline, the stored pose format,
549 the evaluation setup, and the baseline task interface in enough detail to support the main
550 feasibility claims; see Sections 3, 4, and 5.

551 Guidelines:

- 552 • The answer [N/A] means that the paper does not include experiments.
- 553 • If the paper includes experiments, a [No] answer to this question will not be perceived
554 well by the reviewers: Making the paper reproducible is important, regardless of
555 whether the code and data are provided or not.
- 556 • If the contribution is a dataset and/or model, the authors should describe the steps taken
557 to make their results reproducible or verifiable.
- 558 • Depending on the contribution, reproducibility can be accomplished in various ways.
559 For example, if the contribution is a novel architecture, describing the architecture fully
560 might suffice, or if the contribution is a specific model and empirical evaluation, it may
561 be necessary to either make it possible for others to replicate the model with the same
562 dataset, or provide access to the model. In general, releasing code and data is often
563 one good way to accomplish this, but reproducibility can also be provided via detailed
564 instructions for how to replicate the results, access to a hosted model (e.g., in the case
565 of a large language model), releasing of a model checkpoint, or other means that are
566 appropriate to the research performed.
- 567 • While NeurIPS does not require releasing code, the conference does require all submis-
568 sions to provide some reasonable avenue for reproducibility, which may depend on the
569 nature of the contribution. For example
 - 570 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
571 to reproduce that algorithm.
 - 572 (b) If the contribution is primarily a new model architecture, the paper should describe
573 the architecture clearly and fully.
 - 574 (c) If the contribution is a new model (e.g., a large language model), then there should
575 either be a way to access this model for reproducing the results or a way to reproduce
576 the model (e.g., with an open-source dataset or instructions for how to construct
577 the dataset).
 - 578 (d) We recognize that reproducibility may be tricky in some cases, in which case
579 authors are welcome to describe the particular way they provide for reproducibility.
580 In the case of closed-source models, it may be that access to the model is limited in
581 some way (e.g., to registered users), but it should be possible for other researchers
582 to have some path to reproducing or verifying the results.

583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The submission is accompanied by dataset and code resources for the released benchmark and processing pipeline, and these materials are intended to provide the practical entry points needed to reproduce the main experiments and inspect the released asset.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: The paper specifies the task definition, evaluation protocol, data splits, and core baseline settings, including the shared multilingual setup, training schedule, batch size, and learning rate used for the reported results; see Section 5 and Table 1.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: The current experiments report point estimates only; no error bars, confidence intervals, or significance tests are included.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- 635 • The factors of variability that the error bars are capturing should be clearly stated (for
636 example, train/test split, initialization, random drawing of some parameter, or overall
637 run with given experimental conditions).
- 638 • The method for calculating the error bars should be explained (closed form formula,
639 call to a library function, bootstrap, etc.)
- 640 • The assumptions made should be given (e.g., Normally distributed errors).
- 641 • It should be clear whether the error bar is the standard deviation or the standard error
642 of the mean.
- 643 • It is OK to report 1-sigma error bars, but one should state it. The authors should
644 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
645 of Normality of errors is not verified.
- 646 • For asymmetric distributions, the authors should be careful not to show in tables or
647 figures symmetric error bars that would yield results that are out of range (e.g., negative
648 error rates).
- 649 • If error bars are reported in tables or plots, the authors should explain in the text how
650 they were calculated and reference the corresponding figures or tables in the text.

651 8. Experiments compute resources

652 Question: For each experiment, does the paper provide sufficient information on the com-
653 puter resources (type of compute workers, memory, time of execution) needed to reproduce
654 the experiments?

655 Answer: [No]

656 Justification: The paper does not currently report detailed hardware, memory, runtime, or
657 total compute usage for the baseline experiments.

658 Guidelines:

- 659 • The answer [N/A] means that the paper does not include experiments.
- 660 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
661 or cloud provider, including relevant memory and storage.
- 662 • The paper should provide the amount of compute required for each of the individual
663 experimental runs as well as estimate the total compute.
- 664 • The paper should disclose whether the full research project required more compute
665 than the experiments reported in the paper (e.g., preliminary or failed experiments that
666 didn't make it into the paper).

667 9. Code of ethics

668 Question: Does the research conducted in the paper conform, in every respect, with the
669 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

670 Answer: [Yes]

671 Justification: To the best of our knowledge, the work is consistent with the NeurIPS Code of
672 Ethics, and the paper discusses scope, deployment relevance, and limitations of the released
673 resource.

674 Guidelines:

- 675 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
676 Ethics.
- 677 • If the authors answer [No], they should explain the special circumstances that require a
678 deviation from the Code of Ethics.
- 679 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
680 eration due to laws or regulations in their jurisdiction).

681 10. Broader impacts

682 Question: Does the paper discuss both potential positive societal impacts and negative
683 societal impacts of the work performed?

684 Answer: [No]

685 Justification: The paper motivates accessibility and real-world robustness, but it does not yet
686 provide a dedicated broader-impact discussion covering both positive and negative societal
687 effects in a systematic way.

688 Guidelines:

- 689 • The answer [N/A] means that there is no societal impact of the work performed.
- 690 • If the authors answer [N/A] or [No], they should explain why their work has no societal
691 impact or why the paper does not address societal impact.
- 692 • Examples of negative societal impacts include potential malicious or unintended uses
693 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
694 (e.g., deployment of technologies that could make decisions that unfairly impact specific
695 groups), privacy considerations, and security considerations.
- 696 • The conference expects that many papers will be foundational research and not tied
697 to particular applications, let alone deployments. However, if there is a direct path to
698 any negative applications, the authors should point it out. For example, it is legitimate
699 to point out that an improvement in the quality of generative models could be used to
700 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
701 that a generic algorithm for optimizing neural networks could enable people to train
702 models that generate Deepfakes faster.
- 703 • The authors should consider possible harms that could arise when the technology is
704 being used as intended and functioning correctly, harms that could arise when the
705 technology is being used as intended but gives incorrect results, and harms following
706 from (intentional or unintentional) misuse of the technology.
- 707 • If there are negative societal impacts, the authors could also discuss possible mitigation
708 strategies (e.g., gated release of models, providing defenses in addition to attacks,
709 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
710 feedback over time, improving the efficiency and accessibility of ML).

711 11. Safeguards

712 Question: Does the paper describe safeguards that have been put in place for responsible
713 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
714 image generators, or scraped datasets)?

715 Answer: [No]

716 Justification: The paper does not yet include a dedicated safeguards section describing
717 release controls, risk mitigation, or filtering procedures for potentially sensitive scraped web
718 data.

719 Guidelines:

- 720 • The answer [N/A] means that the paper poses no such risks.
- 721 • Released models that have a high risk for misuse or dual-use should be released with
722 necessary safeguards to allow for controlled use of the model, for example by requiring
723 that users adhere to usage guidelines or restrictions to access the model or implementing
724 safety filters.
- 725 • Datasets that have been scraped from the Internet could pose safety risks. The authors
726 should describe how they avoided releasing unsafe images.
- 727 • We recognize that providing effective safeguards is challenging, and many papers do
728 not require this, but we encourage authors to take this into account and make a best
729 faith effort.

730 12. Licenses for existing assets

731 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
732 the paper, properly credited and are the license and terms of use explicitly mentioned and
733 properly respected?

734 Answer: [Yes]

735 Justification: Existing assets are cited in the paper, and the release terms for key reused
736 resources are identified in the accompanying materials. In particular, the Google Research
737 repository hosting YouTube-SL-25 states that datasets are released under CC BY 4.0 and
738 source files under Apache 2.0, while the underlying videos remain subject to the original
739 YouTube content licenses and platform terms.

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

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The paper documents the new pose-native dataset interface, processing pipeline, pose storage schema, and benchmark setting in Sections 3–5, and the released project materials are intended to accompany the asset. The released codebase is distributed under the MIT license, while our released dataset annotations and metadata are distributed under CC BY-NC 4.0 with accompanying usage documentation.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [N/A]

Justification: The paper does not report crowdsourcing studies or research with human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

792 Question: Does the paper describe potential risks incurred by study participants, whether
793 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
794 approvals (or an equivalent approval/review based on the requirements of your country or
795 institution) were obtained?

796 Answer: [N/A]

797 Justification: The paper does not report research with human subjects requiring IRB-style
798 review.

799 Guidelines:

- 800 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
801 with human subjects.
- 802 • Depending on the country in which research is conducted, IRB approval (or equivalent)
803 may be required for any human subjects research. If you obtained IRB approval, you
804 should clearly state this in the paper.
- 805 • We recognize that the procedures for this may vary significantly between institutions
806 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
807 guidelines for their institution.
- 808 • For initial submissions, do not include any information that would break anonymity (if
809 applicable), such as the institution conducting the review.

810 16. Declaration of LLM usage

811 Question: Does the paper describe the usage of LLMs if it is an important, original, or
812 non-standard component of the core methods in this research? Note that if the LLM is used
813 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
814 scientific rigor, or originality of the research, declaration is not required.

815 Answer: [N/A]

816 Justification: The core contribution of the paper is the dataset, processing interface, and
817 baseline evaluation setup rather than an LLM-based method.

818 Guidelines:

- 819 • The answer [N/A] means that the core method development in this research does not
820 involve LLMs as any important, original, or non-standard components.
- 821 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
822 be described.