

The Visual Iconicity Challenge: Evaluating Vision–Language Models on Sign Language Form–Meaning Mapping

Onur Keleş^{1,2}, Aslı Özyürek^{1,3}, Gerardo Ortega⁴, Kadir Gökğöz², Esam Ghaleb^{1,3}

¹Multimodal Language Department, Max Planck Institute for Psycholinguistics

²Department of Linguistics, Boğaziçi University

³Donders Institute for Brain, Cognition and Behaviour, Radboud University

⁴Department of Linguistics and Communication, University of Birmingham

Correspondence: onur.keles1@bogazici.edu.tr and esam.ghaleb@mpi.nl

Abstract

Iconicity, the resemblance between linguistic form and meaning, is pervasive in sign languages, offering a natural testbed for visual grounding in vision–language models (VLMs). We introduce the *Visual Iconicity Challenge*, a video-based benchmark that adapts psycholinguistic measures to evaluate VLMs on three tasks: (i) phonological sign-form prediction, (ii) transparency (inferring meaning from visual form), and (iii) graded iconicity ratings. We assess 17 state-of-the-art VLMs in zero- and few-shot settings on Sign Language of the Netherlands and compare them to human baselines. VLMs mirror human phonological difficulty patterns (e.g., handshape harder than location) and achieve moderate to strong alignment with human iconicity ratings. However, most of them still fail to infer lexical meaning from visual form alone and show a systematic object-based bias that inverts the human preference for action-based signs. Crucially, *models with stronger phonological form prediction correlate better with human iconicity judgments*, indicating shared sensitivity to visually grounded structure. Our findings validate these diagnostic tasks, show that explicit reasoning narrows the open-to-closed-model calibration gap, and motivate human-centric signals for modelling iconicity in multimodal models.

1 Introduction

Language is inherently multimodal: besides speech and text, it includes co-speech gesture and sign languages. Across these modalities, iconicity is the non-arbitrary link between form and meaning. Iconicity can be visual (e.g., speakers use iconic gestures through drawing shapes or trajectories in the air, adding depictive content alongside speech) or even vocal, as in onomatopoeia like “knock knock”, showing that form can transparently reflect meaning (Perlman and Lupyan, 2018). Within sign languages, iconicity is widespread. Estimates

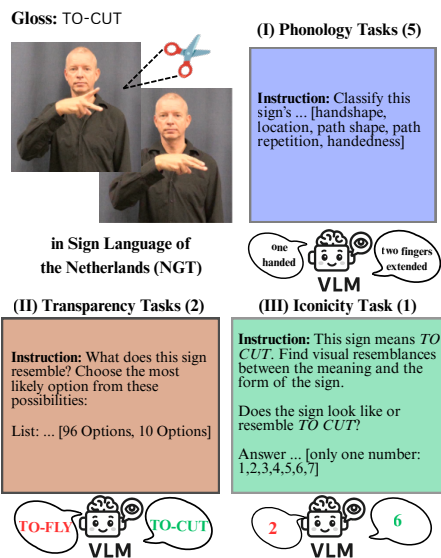


Figure 1: Overview of the *Visual Iconicity Challenge*: evaluation pipeline of the sign TO-CUT in NGT for phonological form prediction, (top right), transparency (bottom left), and iconicity (bottom right) tasks.

suggest that at least a third of lexical signs are iconic (Boyes-Braem, 1986; Campbell et al., 2025) and that between 50–60% of signs’ structure can be directly linked to the physical features of their referents (Ortega, 2017; Pietrandrea, 2002). They depict actions or shapes, providing a natural laboratory for studying the symbol grounding problem: how concepts connect to the physical world (Campbell et al., 2025; Taub, 2001).

For vision–language models (VLMs), sensitivity to form–meaning mapping is a core test of grounding in human-centric signals (Bisk et al., 2020). This is especially relevant for applications in sign language understanding and translation, as well as gesture and action recognition. A capable VLM should attend to *dynamic* bodily movements and hand configurations—not just static objects or text—when interpreting a sign or gesture. sign languages offer a natural testbed for developing and evaluat-

ing models that must perceive temporally extended, simultaneous, visuospatial structure, rather than relying on static context alone (Yin et al., 2021). However, modern VLMs may exhibit static biases, over-attending to appearance features and under-attending to motion cues (Yu et al., 2025), and relying on textual or contextual priors rather than visual evidence when processing gestures (Nishida et al., 2025). To address these questions, we build the **Visual Iconicity Challenge**¹: a sign dataset of Sign Language of the Netherlands (NGT), manually annotated with ground-truth phonological features, iconicity types, and iconicity ratings by non-signers, based on Ortega et al. (2019). The dataset distinguishes between *iconic* signs (with clear visual links to meaning) and *arbitrary* signs (with no visual resemblance).

We evaluate whether models capture different layers of sign–meaning structure, introducing three complementary tasks. Because iconicity links visual form to meaning, it depends on both *phonological form competence* and *analogical reasoning*, requiring models to map structured movement onto conceptual meaning through perceptuo-motor analogy (Thompson and Do, 2019). First, we test whether models can recognise the *phonological form of signs*, including handshape, location, and movement features. Second, we examine *transparency*: if a model can infer a sign’s intended *meaning* from visual form alone (Hoe-mann, 1975), as non-signers often do (Sehyr and Emmorey, 2019). Finally, we test their sensitivity to *iconicity* itself, i.e., whether they can approximate human judgments of graded iconicity. Because previous work shows that VLMs may sometimes rely more on textual or contextual cues than on visual evidence (Nishida et al., 2025), the first two tasks also serve as checks that the models genuinely attend to the visual signal of the sign. Figure 1 shows an overview of the three components. In summary, our contributions are:

- introducing the *Visual Iconicity Challenge*, a benchmark of NGT signs with ground-truth sign phonological annotations and iconicity ratings;
- collecting human baselines for phonology and transparency from a deaf signer and hearing sign-naïve participants;

¹The name is inspired by the “vocal iconicity challenge” of Perlman and Lupyan (2018).

- conducting the first large-scale zero- and few-shot assessment of state-of-the-art VLMs on sign language iconicity, analysing models’ biases for object-based iconicity and failures of form–meaning transparency;
- releasing evaluation code, annotations, and human baselines via a repository for reproducibility and reuse.²

2 Related Work

Iconicity in language and computational models.

Iconicity has long been analysed as structure mapping between form and meaning in sign languages (e.g., depiction of shape or action) (Taub, 2001; Ortega, 2017; Pietrandrea, 2002). Psycholinguistic and lexical studies report substantial iconicity in signed lexicons and roles for iconicity in acquisition, processing, and L2 learning (Boyes-Braem, 1986; Campbell et al., 2025; Karadöller et al., 2024; Caselli and Pyers, 2020).

Recent NLP research has explored analogous patterns in spoken language models. Large language models can capture sound symbolism effects. For example, GPT-4 can generate iconic pseudowords whose meanings humans and models guess above chance (Marklová et al., 2025). Furthermore, larger language models align with human iconicity ratings, indicating some sensitivity to sound symbolism (Loakman et al., 2024). Metaphor understanding, like iconicity, depends on analogical mapping between domains (Lakoff and Johnson, 1980). Tong et al. (2024) introduce the *Metaphor Understanding Challenge*, which tests whether LLMs can interpret metaphors by distinguishing target-domain paraphrases from literal source-domain alternatives. Their findings show that even advanced models often rely on surface similarity rather than analogy.

In the visual modality, sound symbolism studies report weak or dataset-driven effects in CLIP/Stable Diffusion (Alper and Averbuch-Elor, 2023) and mixed evidence for shape/magnitude symbolism in VLMs (Loakman et al., 2024). Extending this understanding to the visual-manual modality of sign languages requires VLMs, which we address in this work.

General multimodal benchmarks. Large-scale multimodal benchmarks have assessed VLM capabilities on image captioning, VQA, and social

²https://github.com/kelesonur/Visual_Iconicity_Challenge

signals. For example, Zhang et al. (2025) introduce MMLA, a suite of 61K multimodal utterances with labels for intent, emotion, style, etc., and report that even fine-tuned state-of-the-art models plateau around 60–70% accuracy. Furthermore, Li et al. (2025) introduce a Multimodal Causal Reasoning benchmark testing whether multimodal models can infer causal relations when crucial evidence appears in visual details. Their results show that models with strong textual reasoning still struggle with visual–conceptual integration. These resources and findings evaluate general multimodal capabilities but do not measure whether models map signed visual form onto meaning or assess graded iconicity relative to human judgments.

Gesture and sign understanding with VLMs.

VLMs underperform on indexical/iconic gestures, especially with visuals-only input, indicating reliance on textual priors (Nishida et al., 2025). Systems like GIRAF mitigate this by injecting structured descriptors (pose skeletons, segmentations, depth) before LLM reasoning, achieving 75% on deictic and 50% on iconic gestures (Lin et al., 2023). Similarly, Zhang et al. (2024) introduce Pose-enhanced VLM, which integrates a skeletal pose modality into a CLIP-like model: one module uses the 2D pose to guide the visual attention to body joints, and another enriches the pose representation with visual context. This integration yields fine-grained action recognition by encouraging the model to focus on human motion cues. In sign language specifically, recent systems fuse additional signals. For example, SignLLM leverages human poses to generate sign language poses for digital human or avatar generation (Fang et al., 2025).

To our knowledge, no prior work has systematically probed VLMs on iconicity in sign languages. Our study is the first to do so at scale, evaluating how well off-the-shelf VLMs perceive the form–meaning transparency that signers exploit.

3 Dataset: The Visual Iconicity Challenge

We present a dataset built on the Sign Language of the Netherlands (NGT) from Karadöller et al. (2024) and Ortega et al. (2019). It contains 96 sign videos (64 iconic signs and 32 arbitrary signs), each with an English gloss (meaning) and human iconicity ratings ranging from 1 to 7 (see the full dataset in Appendix A). This categorisation is based on human iconicity ratings: signs with low ratings ($M = 2.10$, $SD = 0.50$) were classed as arbitrary, and

signs with high ratings ($M = 5.13$, $SD = 1.02$) were classed as iconic.³

Although the benchmark is intentionally compact, the 96 signs were carefully balanced for iconicity, lexical class (verbs vs. nouns), and reference type (object/action/combined/arbitrary), and each sign carries dense psycholinguistic annotations (five phonological parameters, iconicity ratings, and human transparency baselines). This stratified design allows fine-grained per-feature and per-type analyses that would be infeasible at larger scale, and aligns with the size of psycholinguistic stimulus sets in the sign-language literature (Karadöller et al., 2024; Ortega et al., 2019). Our evaluation operationalises iconicity through three complementary tasks targeting different form–meaning mapping levels. (i) **Phonological form prediction** examines whether models perceive the *articulatory structure* of a sign (handshape, location, movement). (ii) **Transparency** asks models to recover a sign’s *lexical meaning* from visual form alone, indexing analogical mapping from form to concept while minimising reliance on linguistic priors. (iii) **Graded iconicity rating** evaluates whether models are sensitive to the *degree* of resemblance between form and meaning by correlating model ratings with human judgments.

Hypothesis. Models that better predict phonological features (e.g., handshape, location, path) should better capture iconicity, since both require grounding in structured bodily properties.

Motivated by the view that iconicity relies on both *phonological form competence* and *analogical reasoning* (Thompson and Do, 2019), we extend the original resource with: (i) detailed phonological annotations for each sign, (ii) iconicity-type labels, and (iii) human baselines for phonology and transparency. These additions support evaluation of VLMs from sub-lexical perception to graded iconicity. See Table 1 for a comparison between the original dataset and our extensions.

³The original stimulus set distinguishes two iconic sub-groups (high vs. low gesture-overlap with normed silent gestures), but as their iconicity ratings do not differ significantly (Karadöller et al., 2024), we treat them as a single iconic category throughout. The categories in the stimulus set were originally determined by consulting a native NGT signer and using a normed silent-gesture database (Ortega and Özyürek, 2020). Also, throughout the paper, M denotes the mean and SD the standard deviation.

Item	Ortega et al.	Ours
Phonology form features		
(based on Klomp and Pfau 2020)		
Handshape	X	✓
Location	X	✓
Path shape	X	✓
Path repetition	X	✓
Handedness	X	✓
Transparency labels (N=96)		
✓	✓	✓
Iconicity		
Ratings (1–7)	✓	✓
Labels (Iconic vs. arbitrary)	✓	✓
Types (e.g., Object or action based)	X	✓
Human baselines		
Phon. form prediction	X	✓
Transparency	X	✓
Iconicity ratings	✓	✓

Table 1: Comparison of the original NGT sign videos dataset (Ortega et al., 2019; Karadöller et al., 2024) and our extensions for the visual iconicity challenge.

3.1 Sign Phonological Form Features

We annotate phonological form features of each sign using a standard NGT phonology framework (Klomp and Pfau, 2020). These are discrete, visual descriptors of articulation (elaboration on the annotation criteria is in Appendix B). In summary, we use five phonological parameters:

- *Handshape*: 7 categories (e.g., fist, flat hand, one finger extended, etc.). Figure 2 illustrates a few categories of the annotated handshapes.
- *Location*: 5 categories of where on the signer’s body or space the sign is articulated, i.e., face/head, torso, arm/shoulder, the opposite hand, or neutral space.
- *Path Shape*: 4 categories of movement trajectory shape, i.e., no movement/hold, straight line, arched curve, and circular motion.
- *Path Repetition*: 2 categories (whether the movement is repeated or only single).
- *Handedness*: 3 categories (one-handed sign, two-handed symmetrical, or two-handed asymmetrical).

A deaf signer and a hearing non-signing researcher performed the annotations. To assess reliability, inter-annotator agreement ranged from 77.9% ($\kappa = 0.73$) for handshape to 98.9% ($\kappa = 0.98$) across parameters. All disagreements were

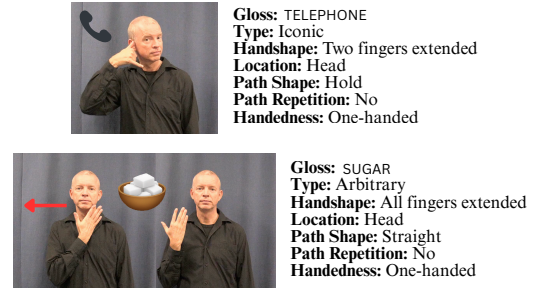


Figure 2: Examples of an iconic vs. an arbitrary sign, with their annotated phonological form features. The sign TELEPHONE is iconic as its form resembles a telephone’s shape, whereas SUGAR is arbitrary with no clear visual link to its meaning.

discussed and resolved. These reliable annotations serve as a gold-standard reference or “ceiling” for assessing how well models can recognise sign form.

Task. Given a sign video, models perform multi-class prediction for each parameter: handshape, location, path shape, path repetition, and handedness (prompts in Appendix C). For this task, we report the accuracy of the model predictions per parameter and the overall average accuracy. This task checks whether the model is capable of extracting form information from the video: where and how signs were articulated.

Human baseline. We gather baseline results from human participants for all the phonological parameters. The 96 stimulus signs were divided into four lists of 24 signs each. Four sign-naïve Turkish undergraduate participants (i.e., without prior knowledge of sign language) were recruited and compensated with course credit, and randomly assigned to lists in a counterbalanced design. Each participant judged all 24 signs in their list on both the phonological feature tasks and the transparency (open-set meaning identification) task. This provided a sign-naïve baseline for both tasks. The human baseline mean phonological accuracy was 0.79 (highest for handedness, lowest for handshape).

3.2 Sign Transparency

Task. Transparency tests whether meaning can be inferred from visual form alone. We use the gloss list (i.e., meaning) of each sign, which is provided in the original dataset. We evaluate two settings: **Transparency₁** (open-set identification among all 96 glosses) and **Transparency₂** (multiple choice with 10 candidates: the target gloss plus 9 distractors). We use accuracy as the primary

metric (proportion of signs correctly identified).

Human baseline. The deaf signer (who annotated the phonological form features and is not a native NGT signer) identified 57/96 glosses in the open-set setting; the sign-naïve group identified 40/96 (same participants and lists as in the phonology baseline). These provide upper- and lower-bound human references for Transparency₁.

3.3 Sign Iconicity Ratings

Task. This task probes whether models capture the degree to which a sign’s form resembles its meaning. We use the original crowdsourced iconicity ratings as **human baselines** (see Appendix A). Each sign has an average iconicity rating on a 1–7 scale (with 7 = “looks exactly like its meaning”, 1 = “not iconic at all”). Models are prompted to produce the iconicity rating for each sign (i.e., the degree of the sign’s resemblance to its meaning). We report Spearman’s ρ (rank alignment) and quadratic-weighted Cohen’s κ_w (ordinal scale agreement) between model ratings and the average human iconicity ratings.

Iconicity types. Iconicity type influences how signers perceive, process, and acquire signs (Ortega et al., 2014, 2017). We annotate each sign for its iconicity type to probe how well models align with these distinctions. These include *object-based* signs ($N = 16$), where the handshape visually resembles a property of the referent (e.g., the wings of a butterfly), and *action-based* signs ($N = 31$), where the hand depicts an action performed on or by the referent (e.g., brushing teeth). The remaining 17 signs belonged to a third category named “combined”, where both strategies were employed for the same sign. The descriptions of these types can be found in Appendix D. These label types enable us to analyse how different iconic strategies affect model predictions and human perception.

4 Models and Inference

Models. We evaluate a representative and diverse set of 11 open-source VLMs and 6 proprietary models, spanning small ($\leq 7B$), medium (27–32B), and large ($\geq 70B$) parameter regimes. Open models: Qwen2.5-VL (72B/32B/7B) (Bai et al., 2025), Qwen3-VL (32B/4B), VideoLLaMA2 (72B/7B) (Cheng et al., 2024), LLaVA-Onevision (72B/7B) (Li et al., 2024), Gemma-3 (27B) (Team et al., 2025b), and MiMo-VL-7B (Team et al., 2025a). Proprietary models: GPT-4o (OpenAI,

2024), GPT-5 (OpenAI, 2025), GPT-5.4 (OpenAI, 2026), Gemini 2.5 Pro (Google DeepMind, 2025a), Gemini 3 Flash (Google DeepMind, 2025b), and Gemini 3.1 Pro (Google DeepMind, 2026). For inference, smaller models ($\leq 7B$) were run on a single NVIDIA A100 GPU, while larger models ($\geq 27B$) were distributed across up to four A100 GPUs. Closed-source models were queried via API calls.

Technical details. Our dataset comprises 96 sign language videos, all recorded at 50 fps, with mean duration of 2.42 ± 0.40 seconds (range: 1.68–3.24 s). All models sample 8 frames per video and process each through their native vision encoder. As a concrete example, Qwen3-VL-32B uses dynamic-resolution encoding, yielding $\sim 1,165$ vision tokens per video (8 frames). In the 4-shot configuration, each shot contains video ($\sim 1,165$ tokens) + question (~ 86 tokens) + answer (~ 1 token) = $\sim 1,252$ tokens, so the total input is 4 shots (5,008 tokens) + test (1,251 tokens) = 6,259 tokens, representing 2.4% of the 262,144-token context window. Other models use similar architectures with varying vision encoders.

Zero-shot setup. All models are evaluated in a zero-shot manner first. We craft a prompt/instruction template for each task that is standardised across models. The prompts explicitly describe the task and the expected answer format, and we ensure the output format is constrained (e.g., just a single number for ratings, or a one-word answer for glosses). For example, for the iconicity rating task, the prompt to the model is:

This sign means: <MEANING>. Some signs are iconic, and some are arbitrary. Find visual resemblances between the meaning and the form of the sign. How much does the sign look like “<MEANING>”? Answer with only one number: 1,2,3,4,5,6,7 (1=not at all, 7=exactly).

We do not use chain-of-thought prompting or specialised prompting tools, as initial trials with those did not show clear benefits. Our aim is to first establish baseline performance; more sophisticated prompting or fine-tuning can be explored in future work. For Qwen3-VL-32B, we additionally test with thinking mode enabled (denoted \dagger) on the iconicity task, to probe whether explicit reasoning improves scale calibration.

Few-shot setup. To examine whether a few examples can improve models’ performance, we

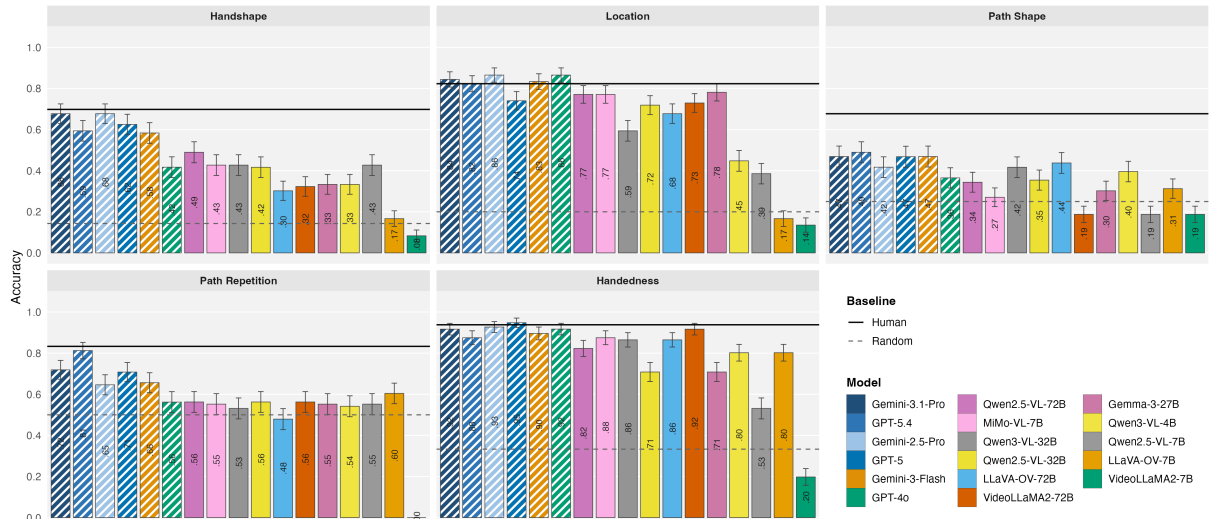


Figure 3: Zero-shot accuracy per phonological feature for all 17 VLMs. Striped bars: closed-source; solid bars: open-source. Solid black line: deaf-native human baseline; dashed grey line: random baseline. Error bars are binomial SEs over 96 signs. Numbers inside bars are mean accuracies.

conduct 4-shot experiments on a representative panel of open models that span the size spectrum: Qwen2.5-VL (72B/32B/7B), Qwen3-VL (32B/4B), Gemma-3-27B, and MiMo-VL-7B. We omit closed models in these settings since few-shot probing suggests that GPT-5 and Gemini variants are already comparatively well-calibrated in zero-shot settings, showing only marginal benefit.

We provide 4 example QA pairs (two iconic signs and two arbitrary signs) before the test query, using the same instruction format and showing the correct outputs for those examples.

5 Results & Discussions

5.1 Phonological Form Prediction

Zero-shot. Figure 3 shows the accuracy of each model on the five phonological sub-tasks. The location of the sign (where on the body the sign is articulated) and handedness (one vs. two hands) were the easiest. For location, 11/17 models reach ≥ 0.70 (best 0.865), while for handedness, 15/17 exceed 0.70. In contrast, handshape and path shape remain the hardest features, with only closed-source models exceeding 0.50 on handshape. Most models exceed the random baseline but remain below the human mean of 0.794. Among closed models, Gemini 3.1 Pro leads with mean accuracy $M = 0.725$, followed by GPT-5.4 ($M = 0.719$), both surpassing Gemini 2.5 Pro ($M = 0.706$); Gemini 3 Flash ($M = 0.688$) falls just below it. Among open-source models, Qwen2.5-VL-72B ($M = 0.598$), MiMo-VL-7B

($M = 0.579$), Qwen3-VL-32B ($M = 0.567$), and LLaVA-OV-72B and Qwen2.5-VL-32B (both $M = 0.552$) lead. Overall, while large models encode phonologically relevant structure, the absolute gap to human performance remains large for most features. Full per-feature results for all models are in Appendix E (Table 8) and Appendix F (Figure 8).

Interestingly, the performance patterns of models for each form feature in Figure 3 mirror well-established acquisitional asymmetries in sign language. Like deaf children and adults (Keleş et al., 2022; Sandler and Lillo-Martin, 2006; Morgan et al., 2007; Marentette and Mayberry, 1999), models find *location* easier than *handshape*.

Few-shot. Few-shot prompting yields modest, model-dependent gains for selected open-source VLMs (Table 2). Qwen2.5-VL-7B shows the largest absolute gain (0.417 \rightarrow 0.550), followed by Qwen2.5-VL-32B and Qwen3-VL-32B, while Qwen2.5-VL-72B changes little. MiMo-VL-7B is essentially flat. This pattern indicates that few-shot prompting mainly benefits models with lower zero-shot performance, presumably by helping them interpret the task format.⁴

⁴A closer look at task-level breakdowns reveals uneven effects: the largest gains occur for path shape and handedness, while location remains unstable, handshape shows minimal improvement, and path repetition is largely unaffected. This suggests that few-shot prompting primarily helps models disambiguate structural features like path shape and handedness.

Model	Mean Accuracy	
	0-shot	4-shot
Qwen2.5-VL-72B	0.598	0.600
Qwen2.5-VL-32B	0.552	0.620
Qwen3-VL-32B	0.567	0.633
Gemma-3-27B	0.535	0.572
MiMo-VL-7B	0.579	0.573
Qwen2.5-VL-7B	0.417	0.550
Qwen3-VL-4B	0.504	0.573

Table 2: Comparison of zero-shot and 4-shot performance on the phonological form prediction task.

5.2 Sign Transparency

Zero-shot. Open-set gloss identification is highly challenging (Table 3). Even the strongest closed-source models perform poorly: the best model, Gemini 3.1 Pro, identifies 28 of the 96 glosses ($\approx 29.2\%$), followed by Gemini 3 Flash (19/96, 19.8%) and Gemini 2.5 Pro (17/96, 17.7%). All remain well below the human baselines (57/96 for the deaf signer of TİD evaluated cross-linguistically, 40/96 for hearing non-signers).⁵ Restricting the task to a 10-way multiple-choice format improves scores (e.g., 47/96 for Gemini 3.1 Pro, 44/96 for Gemini 3 Flash, 42/96 for GPT-5, and 41/96 each for Gemini 2.5 Pro and GPT-5.4). Open-source VLMs perform substantially worse, with the best (Qwen3-VL-32B) achieving only 4/96 correct identifications in the open setting and 19/96 in the 10-way setting. The persistent advantage of closed models over all open-source systems indicates that high-capacity proprietary VLMs are better at leveraging combined visual and linguistic cues.

Across models, correct predictions cluster on visually obvious signs such as TELEPHONE, PISTOL, and TO-WRING (guessed by the large majority of evaluated VLMs), followed by TO-JUGGLE, WIND-SCREEN WIPER, and DEER (see Figure 7 in Appendix F). Interestingly, some arbitrary but cross-linguistically shared signs (e.g., TO-ARGUE, TO-ORDER, PERSON, and TO-DIE) were successfully guessed by a sizable number of VLMs.⁶ Identification of such arbitrary signs suggests that their forms still contain strong visual cues, possibly through conventional metaphorical mappings shared across

⁵Importantly, the deaf signer is a native signer of Turkish Sign Language (TİD) and is *not* familiar with NGT. Their score therefore reflects a Transparency task (inferring meaning from unfamiliar visual forms), not a lexical retrieval task, and serves as a cross-linguistic upper bound on how much iconicity supports inference of meaning across sign languages.

⁶Most of these arbitrary signs were guessed correctly by our human participants too.

sign languages (Meir and Cohen, 2018).

Model	96 options	10 options
Human baselines		
Deaf signer (TİD)	0.594	–
Hearing non-signer	0.417	–
Models		
Gemini-3.1-Pro	0.292	0.490
Gemini-3-Flash	0.198	0.458
GPT-5	0.156	0.438
Gemini-2.5-Pro	0.177	0.427
GPT-5.4	0.146	0.427
GPT-4o	0.073	0.354
Qwen3-VL-32B	0.042	0.198
Qwen2.5-VL-32B	0.052	0.177
MiMo-VL-7B	0.042	0.177
Qwen2.5-VL-72B	0.021	0.167
Qwen3-VL-4B	0.021	0.167
LLaVA-OV-72B	0.031	0.156
VideoLLaMA2-72B	0.031	0.156
Gemma-3-27B	0.021	0.125
VideoLLaMA2-7B	0.010	0.125
Qwen2.5-VL-7B	0.021	0.115
LLaVA-OV-7B	0.021	0.073
Chance (random)	0.010	0.100

Table 3: Transparency task accuracy in 96-option and 10-option conditions. The deaf signer is a native signer of Turkish Sign Language (TİD), *not* NGT.

Few-shot. Four-shot prompting yields no meaningful gains. In the 96-way setting, Qwen2.5-VL (72B & 32B) and Gemma-3-27B each identify 2 of 92 items, and Qwen2.5-VL-7B identifies 1 of 92. In the 10-choice format, the same models score 15–16 of 92, matching their zero-shot levels. These results suggest the bottleneck is not just understanding the task format, but a fundamental limitation in the models’ visual–semantic grounding.

5.3 Sign Iconicity

Zero-shot. For iconicity ratings (Table 4), several models show moderate-to-strong positive correlations with human iconicity judgments ($\rho \geq 0.40$, $p < .001$). Gemini 3.1 Pro achieves the highest rank alignment with human ratings ($\rho = 0.774$), followed by GPT-5.4 ($\rho = 0.612$), GPT-5 ($\rho = 0.607$), and Gemini 3 Flash ($\rho = 0.600$). Among open VLMs, Qwen2.5-VL-72B shows the strongest correlation ($\rho = 0.501$), followed by Qwen3-VL-32B with thinking mode enabled ($\rho = 0.473$), Qwen2.5-VL-7B ($\rho = 0.456$), and Gemma-3-27B ($\rho = 0.452$). Open models such as Qwen2.5-VL-7B ($\kappa_w = 0.421$) and Gemma-3-27B ($\kappa_w = 0.459$) approach mid-tier closed-model calibration (e.g., Gemini 2.5 Pro: $\kappa_w = 0.458$), despite lower

ρ , indicating better scale calibration relative to their rank alignment.

Importantly, model–human alignment is multifaceted: rank consistency (ρ) and ordinal agreement (κ_w) capture complementary aspects of relative ordering and scale calibration. Gemini 3.1 Pro leads on both dimensions ($\rho = 0.774$, $\kappa_w = 0.522$), with GPT-5 and GPT-5.4 following closely ($\kappa_w = 0.511$ and 0.499 , respectively). Several open models (e.g., Qwen2.5-VL-7B, Gemma-3-27B) match closed-model scale calibration (κ_w) despite weaker rank discrimination (ρ), suggesting they use the rating scale appropriately but are less sensitive to relative iconicity differences between signs, while a subset substantially under-rate iconicity (Qwen2.5-VL-32B, VideoLLaMA2-72B, MiMo-VL-7B, VideoLLaMA2-7B). Despite these strengths, most models still compress the scale around the midpoint and systematically over-rate arbitrary signs, thereby reducing the contrast between iconic and arbitrary categories.

Notably, thinking mode substantially benefits Qwen3-VL-32B ($\rho = 0.473$, $\kappa_w = 0.459$ vs. $\rho = 0.356$, $\kappa_w = 0.177$ without thinking), resolving its scale-anchoring failure and lifting it to second place among open models. Its κ_w of 0.459 is joint-highest among open models (tied with Gemma-3-27B) and matches closed-model levels (e.g., Gemini 2.5 Pro: 0.458), suggesting that reasoning helps models adopt a more human-like scale, not just improve relative rankings.

Model	ρ	κ_w
Gemini-3.1-Pro	0.774 ^{***}	0.522
GPT-5.4	0.612 ^{***}	0.499
GPT-5	0.607 ^{***}	0.511
Gemini-3-Flash	0.600 ^{***}	0.396
Gemini-2.5-Pro	0.577 ^{***}	0.458
GPT-4o	0.248 [*]	0.227
Qwen2.5-VL-72B	0.501 ^{***}	0.284
Qwen3-VL-32B [†]	0.473 ^{***}	0.459
Qwen2.5-VL-7B	0.456 ^{***}	0.421
Gemma-3-27B	0.452 ^{***}	0.459
VideoLLaMA2-72B	0.400 ^{***}	0.261
MiMo-VL-7B	0.389 ^{***}	0.179
Qwen3-VL-32B	0.356 ^{***}	0.177
Qwen2.5-VL-32B	0.344 ^{***}	0.152
Qwen3-VL-4B	0.238 [*]	0.122
LLaVA-OV-72B	0.223 [*]	0.172
LLaVA-OV-7B	0.119 ^{ns}	0.112
VideoLLaMA2-7B	0.101 ^{ns}	0.047

Table 4: Graded iconicity rating results. Significance codes for ρ : * $p < .05$, ** $p < .01$, *** $p < .001$, ns $p \geq .05$. [†]Thinking mode enabled during inference.

Few-shot. Few-shot prompting helps Qwen2.5-VL-32B ($\rho : 0.344 \rightarrow 0.510$) and Qwen3-VL-32B ($0.356 \rightarrow 0.488$) most, with smaller but consistent gains for Gemma-3-27B ($0.452 \rightarrow 0.484$) and MiMo-VL-7B ($0.389 \rightarrow 0.412$) (Table 5). Few-shot cues are thus most helpful for models that underperform relative to their capacity.

Model	Spearman ρ	
	0-shot	4-shot
Qwen2.5-VL-72B	0.501	0.521
Qwen3-VL-32B	0.356	0.488
Qwen2.5-VL-32B	0.344	0.510
Gemma-3-27B	0.452	0.484
MiMo-VL-7B	0.389	0.412

Table 5: Comparison of zero-shot and 4-shot performance on the graded iconicity rating.

Type of iconicity. Iconicity is commonly classified by whether a sign depicts an object’s shape or a human action (Ortega et al., 2019). Human raters show a robust preference for *action*-based signs over *object*-based ones, consistent with findings that action signs are acquired earlier and that both deaf children and adults exhibit a cognitive bias toward action-based (handling) iconic forms (Ortega et al., 2017; Sümer and Özyürek, 2025). As illustrated in Figure 4, both humans and large models clearly distinguished arbitrary from iconic signs, indicating that models can broadly recognise iconic structure. However, within iconic signs, differences emerged. Humans show a consistent *action* bias, whereas most open-source models displayed the reverse pattern, favouring *object*-based signs that depict visual features rather than actions. Closed-source models such as Gemini and GPT-5 showed little to no preference between the two types. This static-image bias is well known in computer vision (Zhou et al., 2025), but our contribution is exposing its *linguistic* consequence: a systematic inversion of human iconicity preferences that links a known visual-modeling limitation to a new, linguistically meaningful failure mode (see our qualitative observations in Appendix G).

6 Interaction of Iconicity, Transparency, and Phonology

We hypothesise that models with stronger phonological form predictions are better at both rating graded iconicity and inferring sign meaning, as all three tasks require grounding in structured bodily

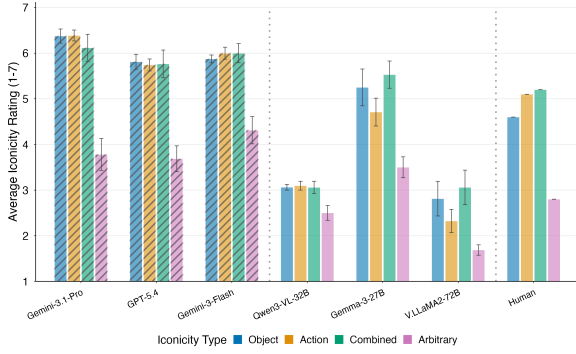


Figure 4: Average iconicity ratings by iconicity type for closed-source (striped) and open-source (solid) models, with the human baseline right of the dotted divider. Error bars: SE across signs. Colours denote reference type: **Object**, **Action**, **Combined**, **Arbitrary**.

properties.

Phonology ↔ Iconicity. As shown in Figure 5, models with higher phonological form accuracies, such as Gemini 3.1 Pro, Gemini 3 Flash, GPT-5.4, GPT-5, and Gemini 2.5 Pro, also achieve closer alignment with human iconicity ratings. Conversely, models with weaker phonological representations (e.g., smaller Qwen and LLaVA variants) show both lower accuracy on phonological features and less consistent treatment of iconicity. Across all 17 models, mean phonological accuracy correlates strongly with iconicity ρ ($\rho = .769$, $p < .001$), suggesting that sensitivity to phonological form and to form–meaning mappings are not independent, but may sign language co-develop (Emmorey, 2014).

Phonology ↔ Transparency. Mean phonological accuracy correlates strongly with transparency accuracy (96-opt: $\rho = .898$; 10-opt: $\rho = .913$; both $p < .001$), indicating that explicit form encoding substantially supports meaning inference regardless of whether the decision space is constrained. Together with the phonology–iconicity link above, our three tasks operationalise complementary aspects of form–meaning mapping: phonological perception of bodily structure, analogical inference from form to concept (transparency), and graded sensitivity to resemblance (iconicity).

This link mirrors the human case, where iconicity is tied to phonological awareness because mapping form features onto conceptual structure requires attending to form and meaning simultaneously. Yet models systematically overrate object-based iconicity (Figure 4), so phonological sensi-

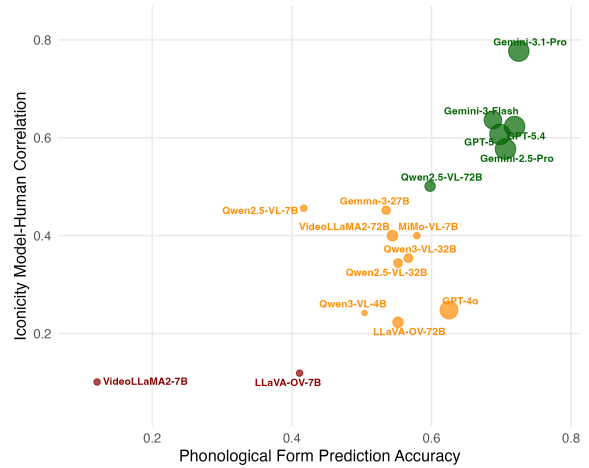


Figure 5: Overall model landscape by zero-shot phonological form prediction accuracy (x -axis) and iconicity Spearman ρ with humans (y -axis). Top-right is best; dot size encodes model size. Colours denote performance tiers: **red** = low, **orange** = mid, **green** = high (assigned by the average of normalised phonology and iconicity scores).

tivity alone does not yield the human action-based preference.

7 Conclusion

We introduced the *Visual Iconicity Challenge*, a diagnostic evaluation that probes phonological form prediction, meaning prediction from form, and iconicity ratings in the Sign Language of the Netherlands. Our evaluations suggest that, compared to human baselines, larger vision large language models mirror human phonological difficulty patterns (e.g., handshape is more difficult than location), can distinguish iconic from arbitrary signs, and correlate moderately to strongly with graded human iconicity ratings. Yet they fail to infer lexical meaning and show a different iconicity-type bias. Bridging this gap requires richer gesture/sign pretraining and dynamic pose encoding; explicit reasoning offers a complementary route, as thinking mode alone lifted open-model scale calibration to closed-model levels.

8 Limitations

Our evaluation has several constraints. The dataset is small (96 isolated NGT signs) with citation-style clips that may not generalise to other sign languages or continuous discourse. Phonological annotations cover five major parameters but omit finer-grained features (orientation, aperture changes, non-manual markers), and the mixed lex-

ical classes (verbs vs. nouns) may affect transparency and iconicity patterns.

Furthermore, we evaluated models only in zero-shot and few-shot settings without fine-tuning, which establishes a diagnostic baseline but likely underestimates potential performance with sign-specific training. Future work may explore fine-tuning, examine model-specific factors (parameter count, memory footprint, mixture-of-experts activation patterns during inference to examine the processing of signs with different levels of iconicity), test robustness under visual perturbations (noise, motion blur), conduct stratified analyses by iconicity type (action-based vs. object-based) and sign difficulty levels, and perform qualitative error analysis to identify whether failures stem from visual perception, analogical reasoning, or lexical-semantic grounding deficits.

Ethics Statement

All human participants provided informed consent prior to participation. The study was approved by the Institutional Review Board in Social Sciences and Humanities at Boğaziçi University.

References

- Morris Alper and Hadar Averbuch-Elor. 2023. [Kiki or bouba? sound symbolism in vision-and-language models](#). 36:78347–78359.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. [Qwen2. 5-vl technical report](#). *arXiv preprint arXiv:2502.13923*.
- Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, and 1 others. 2020. [Experience grounds language](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735.
- Penny Boyes-Braem. 1986. Two aspects of psycholinguistic research: iconicity and temporal structure. In *Proceedings of the Second European Congress on Sign Language Research; Signs of Life*, Amsterdam. University of Amsterdam.
- Erin E. Campbell, Zed Sevcikova Sehyr, Elana Pontecorvo, Ariel Cohen-Goldberg, Karen Emmorey, and Naomi Caselli. 2025. [Iconicity as an organizing principle of the lexicon](#). *Proceedings of the National Academy of Sciences*, 122(16):e2401041122.
- Naomi K Caselli and Jennie E Pyers. 2020. [Degree and not type of iconicity affects sign language vocabulary acquisition](#). *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 46(1):127.
- Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. 2024. [Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms](#).
- Karen Emmorey. 2014. [Iconicity as structure mapping](#). *Philosophical Transactions of the Royal Society B: Biological Sciences*, 369(1651):20130301.
- Sen Fang, Chen Chen, Lei Wang, Ce Zheng, Chunyu Sui, and Yapeng Tian. 2025. [Signllm: Sign language production large language models](#).
- Google DeepMind. 2025a. [Gemini 2.5 pro – model card](#). <https://deepmind.google/models/model-cards/gemini-2-5-pro/>.
- Google DeepMind. 2025b. [Gemini 3 flash – model card](#). <https://deepmind.google/models/gemini/flash/>.
- Google DeepMind. 2026. [Gemini 3.1 pro – model card](#). <https://deepmind.google/models/model-cards/gemini-3-1-pro/>.
- Harry W. Hoemann. 1975. [The transparency of meaning of sign language gestures](#). *Sign Language Studies*, (7):151–161.
- Dilay Z. Karadöller, David Peeters, Francie Manhardt, Asli Özyürek, and Gerardo Ortega. 2024. [Iconicity and gesture jointly facilitate learning of second language signs at first exposure in hearing nonsigners](#). *Language Learning*, 74(4):781–813.
- Onur Keleş, Furkan Atmaca, and Kadir Gökgez. 2022. [Effects of age of acquisition and category size on signed verbal fluency](#). *Language Acquisition*, 29(4):361–383.
- Ulrika Klomp and Roland Pfau, editors. 2020. *A Grammar of Sign Language of the Netherlands (NGT)*, 1 edition. SIGN-HUB Sign Language Grammar Series. SIGN-HUB. Accessed 31-10-2021.
- George Lakoff and Mark Johnson. 1980. [The metaphorical structure of the human conceptual system](#). *Cognitive Science*, 4(2):195–208.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. [Llava-onevision: Easy visual task transfer](#).
- Zhiyuan Li, Heng Wang, Dongnan Liu, Chaoyi Zhang, Ao Ma, JiETING Long, and Weidong Cai. 2025. [Multimodal causal reasoning benchmark: Challenging multimodal large language models to discern causal links across modalities](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 5509–5533, Vienna, Austria. Association for Computational Linguistics.
- Li-Heng Lin, Yuchen Cui, Yilun Hao, Fei Xia, and Dorsa Sadigh. 2023. [Gesture-informed robot assistance via foundation models](#).
- Tyler Loakman, Yucheng Li, and Chenghua Lin. 2024. [With ears to see and eyes to hear: Sound symbolism experiments with multimodal large language models](#). pages 2849–2867.

- Paula F. Marentette and Rachel I. Mayberry. 1999. [Principles for an emerging phonological system: A case study of early asl acquisition](#). In Charlene Chamberlain, Jill P. Morford, and Rachel I. Mayberry, editors, *Language Acquisition by Eye*, pages 71–90. Psychology Press, Mahwah, NJ.
- Anna Marklová, Jiří Milička, Leonid Ryvkin, L'udmila Lacková Bennet, and Libuše Kormaníková. 2025. [Iconicity in large language models](#). *Digital Scholarship in the Humanities*, 40(4):1203–1224.
- Irit Meir and Ariel Cohen. 2018. [Metaphor in sign languages](#). *Frontiers in psychology*, 9:1025.
- Gary Morgan, Sarah Barrett-Jones, and Helen Stoneham. 2007. [The first signs of language: Phonological development in british sign language](#). *Applied Psycholinguistics*, 28(1):3–22.
- Noriki Nishida, Koji Inoue, Hideki Nakayama, Mayumi Bono, and Katsuya Takanashi. 2025. [Do multimodal large language models truly see what we point at? investigating indexical, iconic, and symbolic gesture comprehension](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 514–524, Vienna, Austria. Association for Computational Linguistics.
- OpenAI. 2024. [Gpt-4o system card](https://openai.com/index/gpt-4o-system-card/). <https://openai.com/index/gpt-4o-system-card/>.
- OpenAI. 2025. [Gpt-5 system card](https://cdn.openai.com/gpt-5-system-card.pdf). <https://cdn.openai.com/gpt-5-system-card.pdf>.
- OpenAI. 2026. [Gpt-5.4 thinking system card](https://openai.com/index/gpt-5-4-thinking-system-card). <https://openai.com/index/gpt-5-4-thinking-system-card>.
- Gerardo Ortega. 2017. [Iconicity and sign lexical acquisition: A review](#). *Frontiers in Psychology*, 8:1280.
- Gerardo Ortega and Aslı Özyürek. 2020. [Systematic mappings between semantic categories and types of iconic representations in the manual modality: A normed database of silent gesture](#). *Behavior Research Methods*, 52:51–67.
- Gerardo Ortega, Annika Schiefner, and Aslı Özyürek. 2019. [Hearing non-signers use their gestures to predict iconic form-meaning mappings at first exposure to signs](#). *Cognition*, 191:103996.
- Gerardo Ortega, Beyza Sumer, and Aslı Özyürek. 2014. [Type of iconicity matters: Bias for action-based signs in sign language acquisition](#). In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 36.
- Gerardo Ortega, Beyza Sümer, and Aslı Özyürek. 2017. [Type of iconicity matters in the vocabulary development of signing children](#). *Developmental Psychology*, 53(1):89.
- Marcus Perlman and Gary Lupyan. 2018. [People can create iconic vocalizations to communicate various meanings to naïve listeners](#). *Scientific Reports*, 8(1):2634.
- Paola Pietrandrea. 2002. [Iconicity and arbitrariness in Italian Sign Language](#). *Sign Language Studies*, 2(3):296–321.
- Wendy Sandler and Diane Carolyn Lillo-Martin. 2006. *Sign language and linguistic universals*. Cambridge University Press.
- Zed Sevcikova Sehyr and Karen Emmorey. 2019. [The perceived mapping between form and meaning in American Sign Language depends on linguistic knowledge and task: Evidence from iconicity and transparency judgments](#). *Language and Cognition*, 11(2):208–234.
- Beyza Sümer and Aslı Özyürek. 2025. [Action bias in describing object locations by signing children](#). *Sign Language & Linguistics*.
- Sarah F Taub. 2001. *Language from the body: Iconicity and metaphor in American Sign Language*. Cambridge University Press.
- Core Team, Zihao Yue, Zhenru Lin, Yifan Song, Weikun Wang, Shuhuai Ren, Shuhao Gu, Shicheng Li, Peidian Li, Liang Zhao, Lei Li, Kainan Bao, Hao Tian, Hailin Zhang, Gang Wang, Dawei Zhu, Cici, Chenhong He, Bowen Ye, and 55 others. 2025a. [Mimo-v1 technical report](#). *Preprint*, arXiv:2506.03569.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025b. [Gemma 3 technical report](#). *Preprint*, arXiv:2503.19786.
- Arthur Lewis Thompson and Youngah Do. 2019. [Defining iconicity: An articulation-based methodology for explaining the phonological structure of ideophones](#). *Glossa: a journal of general linguistics*.
- Xiaoyu Tong, Rochelle Choenni, Martha Lewis, and Ekaterina Shutova. 2024. [Metaphor understanding challenge dataset for LLMs](#). *Preprint*, arXiv:2403.11810.
- Kayo Yin, Amit Moryossef, Julie Hochgesang, Yoav Goldberg, and Malihe Alikhani. 2021. [Including signed languages in natural language processing](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7347–7360, Online. Association for Computational Linguistics.
- Yating Yu, Congqi Cao, Yifan Zhang, and Yanning Zhang. 2025. [Learning to generalize without bias for open-vocabulary action recognition](#). *Preprint*, arXiv:2502.20158.
- Hanlei Zhang, Zhuohang Li, Yeshuang Zhu, Hua Xu, Peiwu Wang, Haige Zhu, Jie Zhou, and Jinchao Zhang. 2025. [Can large language models help multimodal language analysis? mmla: A comprehensive benchmark](#).
- Haosong Zhang, Mei Chee Leong, Liyuan Li, and Weisi Lin. 2024. [Pevl: Pose-enhanced vision-language model for fine-grained human action recognition](#). In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18857–18867.
- Shijie Zhou, Alexander Vilesov, Xuehai He, Ziyu Wan, Shuwang Zhang, Aditya Nagachandra, Di Chang, Dongdong Chen, Eric Xin Wang, and Achuta Kadambi. 2025. [Vlm4d: Towards spatiotemporal awareness in vision language models](#). In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8600–8612.

A Stimuli

Iconic Signs (n = 64)

Sign	Iconicity rating	Sign	Iconicity rating
TO-BREAK	6.79	TO-INJECT	5.0
TO-CRY	6.74	TO-STAPLE	4.89
WINDSCREEN WIPER	6.63	CALCULATOR	4.8
TO-CUT	6.61	PENGUIN	4.78
ELEPHANT	6.53	RATTLE	4.75
BICYCLE	6.44	CAR	4.7
BIRD	6.42	CURTAINS	4.7
BABY	6.39	BRIDGE	4.6
KEY	6.26	DEER	4.6
TELEPHONE	6.22	HELICOPTER	4.6
TO-WRING	6.12	MONKEY	4.5
TO-SWIM	6.11	SPIDER	4.5
TO-SLAP	6.11	ZIMMER	4.42
TO-PUMP	6.11	TO-ERASE	4.2
PIANO	6.05	TO-SMS	4.2
TO-KNOCK	6.05	PLANE	4.11
BUTTERFLY	5.94	BALL	4.05
TO-CRASH	5.79	DOOR	4.0
SNAKE	5.74	WHEELCHAIR	4.0
TO-FLY	5.74	CHICKEN	3.83
TABLE	5.7	BLANKET	3.8
PISTOL	5.61	CELL	3.8
EAGLE	5.53	DRILL	3.8
TO-CUT	5.51	TO-PLAY-CARDS	3.8
LAPTOP	5.44	BOTTLE	3.68
UMBRELLA	5.42	CAT	3.61
TO-JUGGLE	5.42	SUITCASE	3.6
CAMEL	5.4	LOBSTER	3.5
SPOON	5.3	TO-PUT-CLOTHES-ON	3.32
TO-STEAL	5.11	BED	3.21
TOWEL	5.1	RESTAURANT	3.21
BOX	5.05	RABBIT	3.16

Arbitrary Signs (n = 32)

Sign	Iconicity rating	Sign	Iconicity rating
AMBULANCE	3.11	MUMMY	2.06
TO-ARGUE	2.94	KIWI	2.06
BEAR	2.89	TO-GOSSIP	2.0
TO-SHOUT	2.79	TO-GO-OUT	1.89
INTERPRETER	2.74	TOILET	1.79
DOG	2.69	ELECTRICITY	1.79
TO-DIE	2.58	PRAM	1.74
PERSON	2.53	DOCTOR	1.74
TO-ORDER	2.44	BUS	1.74
TREE	2.26	HORSE	1.67
TO-LAUGH	2.26	WATER	1.63
SOFA	2.26	BUILDING	1.63
ROOM	2.22	PUPPET	1.53
SHEEP	2.16	FRUIT	1.47
FIRE	2.11	SUGAR	1.37
TO-COOK	2.11	LIGHTBULB	1.22

B Criteria for Phonological Feature Annotation

The following guidelines summarize the decision criteria we applied when annotating the five phonological features of each NGT sign. Our annotations were mainly based on the descriptions drawn from the phonology chapters of *A Grammar of Sign Language of the Netherlands (NGT)* (Klomp and Pfau, 2020). We follow the general phonological descriptions in the NGT grammar but use our own simplified label set for annotation and model evaluation.

Handshape: Handshapes were coded using seven discrete labels:

- All fingers closed to a fist
- All fingers extended
- All fingers curved or clawed
- One (selected) finger extended
- One (selected) finger curved or clawed
- Two or more (selected) fingers extended
- Two or more (selected) fingers curved or clawed

These categories are drawn from the NGT phonological inventory, but we simplify them by collapsing sub-types and by omitting features such as orientation or aperture change. However, our labels treat each sign as having a single static handshape; they therefore do not fully capture signs in which the handshape itself changes over time. For example, signs where a fist closes or opens during the articulation. For such dynamically changing signs we accepted *multiple answers as correct*, so that both start and end configurations are treated as valid.

Location: Each sign was assigned to one of five major location categories:

- Hands touching head/face
- Hands touching torso
- Hands touching arm
- Hands touching weak/passive hand
- Hands in front of the body or face (neutral space)

If a sign involved contact with multiple regions, the primary lexical target location was coded.

Path Shape: Primary path movement was classified using four labels:

- Hold: no path or directional movement
- Straight: linear horizontal, vertical, or diagonal trajectory
- Arched: curved or semicircular trajectory
- Circular: full or near-full circular path

Path Repetition: Repetition of the movement was coded as:

- Single: one primary stroke
- Repeated: movement is repeated

Handedness: Handedness was coded according to the two-handed typology:

- One-handed
- Two-handed symmetrical: both hands share the same handshape and movement
- Two-handed asymmetrical: hands differ in handshape and/or movement

C Used Prompts

Phonological Form Prediction Instructions:

Handshape? H1=all fingers closed to a fist, H2=all fingers extended, H3=all fingers curved or clawed, H4=one (selected) finger extended, H5=one (selected) finger curved or clawed, H6=two or more (selected) fingers extended, H7=two or more(selected) fingers curved or clawed)

Location? Major sign location? Answer with only one: L1, L2, L3, L4, L5 (L1=hands touching head/face, L2=hands touching torso, L3=hands touching arm, L4=hands touching weak/passive hand, L5=hands in front of the body or face)

Path Shape? Movement path shape? Answer with only one: Hold, Straight, Arched, Circular. (Hold=no path or direction, Straight=move in a straight line, Arched=move in an arched line, Circular=move in a circular path)

Path Repetition? Answer with only one: Single, Repeated. (Single=one movement, Repeated=multiple or repeated movements)

Handedness? Answer with only one: One-handed, Two-handed symmetrical, Two-handed asymmetrical. (One-handed=only one hand is used in the sign, Two-handed symmetrical=two hands are used but the hands move together and have the same handshape, Two-handed asymmetrical=two hands are visible, but one hand does not move and the hands have different handshapes)

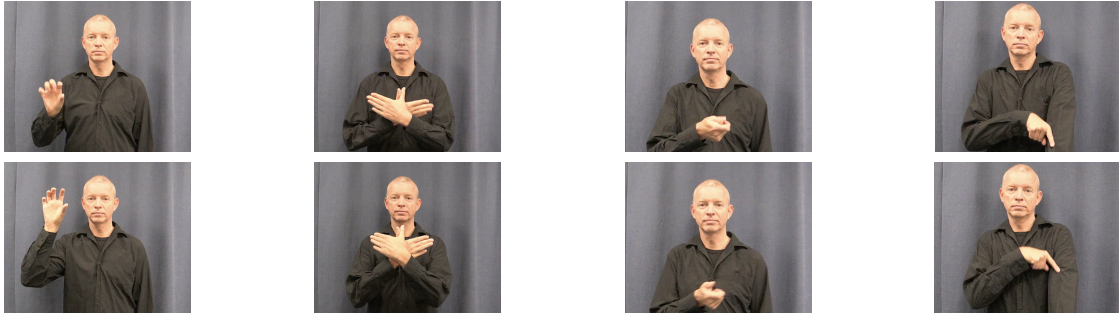
Transparency (open-set over 96 glosses) and Transparency₂ (10-choice Instructions)

What does this sign resemble? Look at the form and movement of the sign. Choose the most likely option from these possibilities: <OPTIONS>. Answer with only the exact word from the list.

Iconicity Rating Instructions:

This sign means: <MEANING>. Some signs are iconic and some are arbitrary. Find visual resemblances between the meaning and the form of the sign. How much does the sign look like “<MEANING>”? Answer with only one number: 1,2,3,4,5,6,7 (1=not at all, 7=exactly).

D Iconicity Type Examples



Combined iconic

SPIDER: The wiggling motion of the hands conveys the spider’s movement, while the curved fingers depict its legs.

Object-based iconic

BUTTERFLY: two hands mirror the referent’s wings.

Action-based iconic

TO-SMS: thumb and fingers enact the *typing/texting* action.

Arbitrary

ELECTRICITY: hand configuration and path show no transparent visual resemblance; the form–meaning link is purely conventional.

Figure 6: Representative frames illustrating the four iconicity categories. Each pair of frames shows how the sign fits its category.

E Full Results Tables for Phonology

Model	Handshape	Location	Path Shape	Path Repetition	Handedness	Mean
Human baseline (hearing non-expert)	0.698	0.823	0.677	0.833	0.938	0.794
<i>Closed-source models</i>						
Gemini-3.1-Pro	0.677	0.844	0.469	0.719	0.917	0.725
GPT-5.4	0.594	0.823	0.490	0.812	0.875	0.719
Gemini-3-Flash	0.583	0.833	0.469	0.656	0.896	0.688
Gemini-2.5-Pro	0.677	0.865	0.417	0.646	0.927	0.706
GPT-5	0.625	0.740	0.468	0.708	0.948	0.698
GPT-4o	0.417	0.865	0.365	0.562	0.917	0.625
<i>Open-source models</i>						
Qwen2.5-VL-72B	0.490	0.771	0.344	0.563	0.823	0.598
MiMo-VL-7B	0.427	0.771	0.271	0.552	0.875	0.579
Qwen3-VL-32B	0.427	0.594	0.417	0.531	0.865	0.567
LLaVA-OV-72B	0.302	0.677	0.438	0.479	0.865	0.552
Qwen2.5-VL-32B	0.417	0.719	0.354	0.563	0.708	0.552
VideoLLaMA2-72B	0.323	0.729	0.188	0.563	0.917	0.544
Gemma-3-27B	0.333	0.781	0.302	0.552	0.708	0.535
Qwen3-VL-4B	0.333	0.448	0.396	0.542	0.802	0.504
Qwen2.5-VL-7B	0.427	0.385	0.188	0.552	0.531	0.417
LLaVA-OV-7B	0.167	0.167	0.313	0.604	0.802	0.411
VideoLLaMA2-7B	0.083	0.135	0.188	0.000	0.198	0.121
Random baseline	0.143	0.200	0.250	0.500	0.333	0.285

Table 8: Phonological form prediction accuracy by model and phonological subtasks, with random and human baselines, and mean accuracy across all subtasks.

F Additional Figures

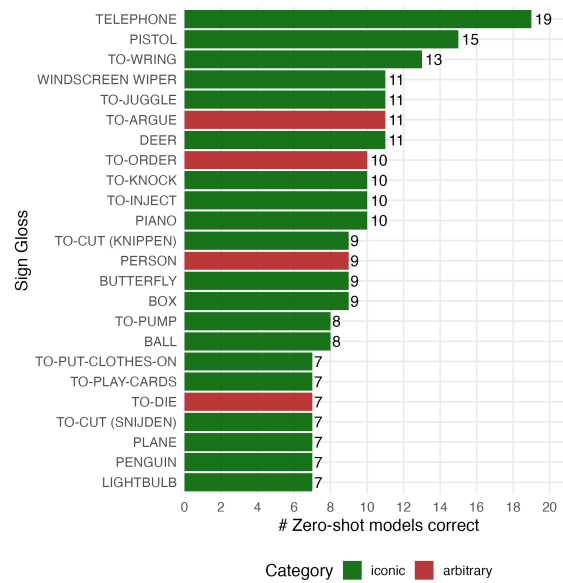


Figure 7: Correctly guessed signs from sign video only (≥ 7 VLMs) in the Transparency₂ Task. Bars in green mark iconic signs; red bars mark arbitrary signs that were nonetheless guessed by multiple VLMs.



Figure 8: Zero-shot accuracy of the top six performing models (3 closed: Gemini 3.1 Pro, Gemini 3 Flash, GPT-5.4; 3 open: Qwen2.5-VL-72B, MiMo-VL-7B, Qwen3-VL-32B) across five phonological features, averaged over 96 signs. Solid colored lines are models; the dashed black line is the human (deaf-native) baseline, and the grey dotted line is the random baseline.

G Qualitative Error Analysis

To complement the quantitative results, we conducted a qualitative inspection of signs with the highest discrepancy between human and the best-performing open-source VLM ratings. We identified three recurring failure modes which we illustrate with representative examples in Table 9. The categories highlight *what* models fail to perceive in signed forms, not only *how often* they fail.

Missing Agent	NGT signs often depict an agent manipulating an object through transitive actions. VLMs appear to be trained predominantly on static, isolated objects and fail to map the holding or manipulating action to the referent noun. <i>Examples:</i> BABY (cradling motion); CALCULATOR (finger pressing keys on held object); TOWEL (wiping motion); CAR (gripping and turning steering wheel).
Static Bias	Models rate signs with static handshapes higher than those with dynamic motions, suggesting better alignment with static visual contours than dynamic event simulations. VLMs privilege shape-matching over motion-based semantic grounding, inverting the human preference for action-based iconicity. <i>Example:</i> TO-FLY (static airplane classifier handshape rated higher than dynamic BIRD with flapping wings).
Gloss Sensitivity	Visual forms that are historically specific, culturally bound, or incongruent with contemporary web-scraped training data cause systematic failures. <i>Examples:</i> TO-SMS (single-handed typing on old keypad phone; web data shows two-thumbed smartphone typing); ZIMMER/ROLLATOR (walking frame uncommon in training data); TO-STEAL (pickpocketing with peripheral twisting motion; web images show generic theft scenes).

Table 9: Three recurring qualitative failure modes observed in zero-shot VLM responses on signs with the largest model–human discrepancy. The *Missing Agent* and *Static Bias* modes correspond to the under-rating of action-based iconicity in Figure 4; *Gloss Sensitivity* reflects limitations of web-scraped training data.