



PointArena: Probing Multimodal Grounding Through Language-Guided Pointing

Long Cheng^{1*} Jiafei Duan^{1,2*} Yi Ru Wang^{1†} Haoquan Fang^{1,2†} Boyang Li^{1†}
Yushan Huang¹ Elvis Wang³ Ainaz Eftekhar^{1,2} Jason Lee^{1,2} Wentao Yuan¹
Rose Hendrix² Noah A. Smith^{1,2} Fei Xia¹ Dieter Fox¹ Ranjay Krishna^{1,2}

¹University of Washington ²Allen Institute for Artificial Intelligence

³Anderson Collegiate Vocational Institute

<https://pointarena.github.io>

Abstract

Pointing serves as a fundamental and intuitive mechanism for grounding language within visual contexts, with applications spanning robotics, assistive technologies, and interactive AI systems. While recent multimodal models have begun supporting pointing capabilities, existing benchmarks typically focus only on referential object localization. We introduce **PointArena**, a comprehensive platform for evaluating multimodal pointing across diverse reasoning scenarios. PointArena comprises three components: (1) **Point-Bench**, a curated dataset of approximately 1,000 pointing tasks across five reasoning categories; (2) **Point-Battle**, an interactive web-based arena facilitating blind, pairwise model comparisons, which has collected over 4,500 anonymized votes; and (3) **Point-Act**, a real-world robotic manipulation system allowing users to directly evaluate model pointing in practical settings. We conducted extensive evaluations of both state-of-the-art open-source and proprietary models. Results indicate that Molmo-72B consistently outperforms others, though proprietary models increasingly demonstrate comparable performance. Additionally, we find that supervised training targeting pointing tasks significantly improves performance. Across our multi-stage evaluation pipeline, we observe strong correlations, underscoring the critical role of precise pointing in enabling multimodal models to bridge abstract reasoning with real-world actions.

1 Introduction

Pointing focuses our attention. It is one of the earliest and most universal non-verbal methods we use to communicate intent; in fact, children learn to point as a prelinguistic form of communication [29]. Precise spatial grounding—*pointing*—enables a wide range of practical and high-impact applications across robotics, assistive technology, human-computer interaction, and vision-language interfaces. In robotics, a pointing-capable model can interpret language commands like “pick up the red cup next to the bowl” and translate them into precise spatial actions [37], enabling fine-grained object manipulation in cluttered environments [13]. In assistive technologies, systems can help visually impaired users by answering spatial queries such as “where is the handle on this door?” [6] or ‘which one is the garlic?’ In education or creative tools, pointing allows for interactive visual tutoring, such as identifying components in scientific diagrams or guiding a learner through a painting [16]. Even in everyday virtual assistants or search engines, the ability to refer to specific image regions via pointing could make multimodal interactions more intuitive and expressive [10]. Across these

*Co-first authors.

†Co-second authors.

32 domains, pointing provides a low-bandwidth yet powerful spatial interface for grounding language in
 33 vision—precise enough for manipulation, intuitive enough for communication, and general enough
 34 to scale with modern multimodal models.

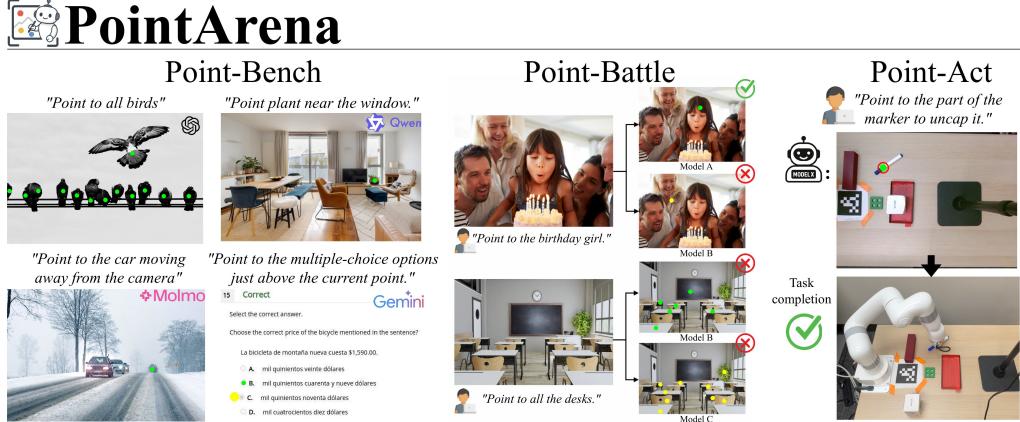


Figure 1: **Overview of PointArena.** PointArena consists of three components: **Point-Bench**, a curated dataset for evaluating grounded pointing across five reasoning types; **Point-Battle**, a live platform for blind, pairwise model comparisons with user voting; and **Point-Act**, real-world task involving manipulation via pointing-based language commands.

35 Recent advances in multimodal models have begun to incorporate more dynamic and spatially
 36 expressive forms of interaction. The Segment Anything Model (SAM) [19] enables segmentation from
 37 sparse visual prompts such as points or boxes, revealing the potential of fine-grained spatial control.
 38 Google’s Gemini models [14] push the boundaries of long-context visual reasoning, incorporating
 39 multiple modalities over extended sequences. In parallel, new datasets have emerged to support
 40 explicit spatial referencing. Molmo’s PixMo dataset [10] introduces 2D pointing as a form of
 41 multimodal alignment between images and instructions, while RoboPoint [37] focuses on spatial
 42 affordance prediction by linking instructions to interaction-relevant keypoints in robotic contexts.
 43 These setups bias evaluations toward pixel-level accuracy rather than conceptual reasoning and often
 44 lack diversity or scalability.

45 There is a need for a holistic evaluation platform to make progress towards language-guided pointing.
 46 Although datasets for referring expressions exist (e.g. RefCOCO, RefCOCO+, and RefCOCOg [17,
 47 34]), they are focused on a subset of pointing tasks: object location. They lack the ambiguity and
 48 contextual variability that users expect from modern interactive models, limiting their utility for
 49 studying pragmatic or interactive applications. As a result, they offer only partial insight into the full
 50 spectrum of grounding required for embodied or assistive agents.

51 We, therefore, propose **PointArena**, a platform to probe and evaluate grounded visual reasoning with
 52 pointing. PointArena presents a suite of tasks where a multimodal model must answer questions
 53 or resolve instructions by combining language and pointing gestures to identify specific image
 54 regions. These tasks go beyond traditional VQA by requiring spatial outputs (e.g., selecting a
 55 location or region) rather than purely textual ones. PointArena allows for both unambiguous and
 56 ambiguous scenarios, supporting studies of disambiguation, spatial commonsense, and pragmatic
 57 inference. Unlike bounding boxes, segmentation masks, or free-form text responses, pointing offers
 58 high-precision signal that avoids reliance on object contours or dense annotations and is directly
 59 compatible with human evaluation.

60 **PointArena** decomposes pointing into three stages of evaluation: 1) **Point-Bench** is a curated
 61 dataset of 982 manually selected, annotated, and verified image-question pairs across five high-level
 62 categories (Spatial, Affordance, Counting, Steerable, and Reasoning), an interactive, online platform
 63 for blind, pairwise comparison between models based on user instructions. Users select from curated
 64 or custom-uploaded images. Voting is anonymized, and we have collected over 4,500 votes from
 65 more than 100 participants. 3) **Point-Act** is a real-world benchmark that evaluates the utility of
 66 pointing in for a downstream application. The system directs a robotic arm to manipulate objects

67 through pointing-based language commands. All three evaluation stages require minimal human
68 effort; each is self-contained and can run live to evaluate any model.

69 Through our evaluation of both open-source and proprietary models across the three stages of the
70 PointArena benchmark, we find that Molmo-72B achieves the highest performance on Point-Bench,
71 with proprietary models such as Gemini-2.5-Pro performing comparably. Models trained with
72 explicit pointing supervision consistently outperform those without. We also observe a strong
73 correlation between static benchmark accuracy and human preference in Point-Battle. Notably, we
74 find that adding language reasoning (e.g., Chain-of-Thought [31]) does not improve visual grounding
75 for pointing tasks. Our study further reveals several other actionable insights into model behavior
76 and evaluation design. We see **PointArena** as a missing component necessary as we develop
77 general-purpose vision-language models that can reason about and interact with the world.

78 2 Related work

79 **Grounding benchmarks.** There are many benchmarks for visual grounding and spatial reasoning
80 capabilities of multimodal large language models (MLLM). The RefCOCO, RefCOCO+, and Ref-
81 COCOg datasets focus on 2D visual grounding, with RefCOCOg emphasizing long-form referring
82 expressions and fine-grained object distinctions [35]. In 3D, ScanRefer provides 51,583 descriptions
83 across 800 RGB-D scans, supporting joint language-geometry models for 3D object localization [7].
84 ReferIt3D and CityRefer extend this to fine-grained and outdoor settings, with CityRefer incorporat-
85 ing geographic features from OpenStreetMap [2, 23]. Interactive benchmarks such as GuessWhat?!
86 evaluate multi-turn object grounding through binary dialog across 150K games [9]. Flickr30K Enti-
87 ties supports phrase-region grounding through 276K bounding boxes and cross-caption coreference
88 annotations [26]. These datasets primarily focus on bounding box grounding, object retrieval, or
89 dialog-based localization. They do not explicitly evaluate pointing behavior.

90 **Arena style evaluation.** Arena-style evaluations have emerged as a popular method for comparing
91 large language models (LLMs) through anonymized pairwise comparisons and user voting, ini-
92 tially popularized by *Chatbot Arena*, which amassed over 240K votes across diverse models and
93 languages [8]. Extensions like *MT-Bench* added multi-turn dialogues with AI-driven scoring [39],
94 while *am-ELO* stabilized rankings through maximum likelihood estimation and annotator reliability
95 modeling [21]. Fully automated frameworks such as *Auto-Arena*, leveraging LLM-generated content
96 and voting, closely mirrored human judgments [38]. *BenchBuilder* automated benchmark creation
97 from crowdsourced data, enhancing statistical robustness. Specialized variants like *Werewolf Arena*
98 examined social reasoning [5], and *OpenArena* facilitated offline evaluations [27]. Despite their
99 scalability and user alignment, arena-style methods lack ground-truth validation and are vulnera-
100 ble to manipulation [22], prompting alternative approaches such as *Tournament Evaluation* [18].
101 *PointArena* uniquely addresses these issues by combining ground-truth evaluation (Point-Bench) with
102 human-preference arena comparisons (Point-Battle).

103 **Models that point or sketch.** Multimodal large language models (MLLMs) have made notable
104 progress in vision-language reasoning, with GPT-4V showcasing tasks like story generation and
105 OCR-free math [33]. Most use encoder-decoder architectures with visual-language fusion to enhance
106 cross-modal integration. MiniGPT-4 aligns ViT features with language models via Q-Formers [30],
107 while Molmo improves 2D spatial grounding by regressing normalized coordinates, reaching 92%
108 precision in icon localization [11]. RoboPoint [36] instruction-tunes VLMs on synthetic data for
109 robotic affordance prediction, outperforming GPT-4o and PIVOT [24] by over 20% in spatial accuracy.
110 Other models, like VisCPM and Qwen-VL, add region-level control and multilingual support, while
111 NEt-GPT integrates 3D, audio, and video inputs [33]. Yet, challenges in precise spatial localization
112 persist, raising questions about what drives effective pointing in MLLMs.

113 3 PointArena

114 Evaluating the ability of MLLMs to localize language-referred entities in images requires benchmarks
115 that are both precise and diagnostic. Existing benchmarks often emphasize classification or captioning,
116 but fall short when it comes to assessing fine-grained spatial grounding—the ability to resolve
117 natural language instructions into specific image coordinates. This capability is critical not only for
118 understanding model alignment with human intent, but also for enabling downstream applications in

119 robotics [36], augmented reality [12], and interactive web agents [15], and potentially contributing
 120 to explainability [25]. As both specialized pointing models and general-purpose MLLMs improve,
 121 standardized evaluation across open-source and proprietary systems becomes essential.

122 We introduce **PointArena**, an evaluation suite for language-conditioned pointing, comprising three
 123 stages: (i) **Point-Bench**, a curated dataset for controlled measurement of spatial localization accuracy;
 124 (ii) **Point-Battle**, a live, blinded human-preference arena for pairwise model comparison; and (iii)
 125 **Point-Act**, a real-world robotic setting that evaluates pointing precision through physical execution.
 126 Together, these components form a unified framework to quantify and analyze how well MLLMs
 127 ground language into visual and physical space.

128 3.1 Task formulation

129 We formalize pointing as a language-conditioned fine-grained localization task. The input consists
 130 of an RGB image $I \in \mathbb{R}^{H \times W \times 3}$ and a natural-language instruction prompt $q = \{w_t\}_{t=1}^T$. A
 131 multimodal large language model (MLLM) \mathcal{F}_θ takes (I, q) as input and predicts a set of image-space
 132 coordinate points $P = \{(x_i, y_i)\}_{i=1}^K$, where each point lies within the image bounds: $x_i \in [0, W-1]$,
 133 $y_i \in [0, H-1]$.

134 Ground-truth supervision is provided as a set of binary masks $\{M_j\}_{j=1}^{K^*}$, with each mask $M_j \in$
 135 $\{0, 1\}^{H \times W}$ denoting the valid region for one of K^* annotated targets. A predicted point (x_i, y_i) is
 136 considered correct if it lies within the spatial support of some mask M_j , i.e., $M_j[y_i, x_i] = 1$.

137 A prediction is considered *successful* if:

- 138 1. The number of predicted points matches the number of target regions: $K = K^*$,
- 139 2. Each target region M_j is covered by at least one predicted point: $\exists(x_i, y_i) \in P$ such that $M_j[y_i, x_i] = 1$.

141 This formulation enables fully automated evaluation given access to the ground-truth masks, with no
 142 need for human annotators at test time.



Figure 2: **Overview of the five Point-Bench categories and the annotation UI.** Point-Bench includes 982 image-query pairs grouped into five categories: **Spatial** (positional references), **Affordance** (functional part identification), **Counting** (attribute-based grouping), **Steerable** (relative pointing), and **Reasoning** (open-ended visual inference). Each example shows a representative query and the corresponding target. On the right, we show the Gradio-based annotation interface used to collect and refine segmentation masks. Initial masks are generated using SAM and refined by annotators, followed by manual verification.

143 3.2 Point-Bench

144 Point-Bench is the largest benchmark for evaluating language-guided pointing, comprising 982
 145 text-image pairs with pixel-level target masks collected from public sources after April 20, 2025.
 146 The dataset is evenly divided into five task-driven categories—Spatial, Affordance, Counting, Steer-
 147 able, and Reasoning—derived from a survey of question types frequently tackled by open-source
 148 MLLMs [11, 36, 28]. Each category targets a distinct capability: 1) *Spatial* focuses on positional

149 queries within scenes rich in spatial relationships or repeated objects (e.g., “Point to the leftmost tree
150 in the image”); 2) *Affordance* emphasizes functional parts of objects, typically in tabletop scenes,
151 prompting queries like “Point to the handle used for pouring”; 3) *Counting* features multiple similar
152 items and supports queries about subsets based on number or attributes, such as “Point to all the blue
153 cars in the image”; 4) *Steerable* leverages images from the PixMo dataset that include a reference
154 point, guiding annotators to ask relative-position questions like “Point to the item closest to the
155 marked point”; and 5) *Reasoning* presents event-rich or abstract scenes, inviting open-ended queries
156 that require inference, such as “Point to the tallest man-made object in the image.” Annotators,
157 recruited via crowdsourcing, were free to ask any question, but carefully selected category-specific
158 images naturally guided them toward prompts aligned with each reasoning axis. These curated splits
159 together support systematic evaluation of an MLLM’s ability to recognize, reason, and precisely
160 ground language in visual space.

161 To construct a Point-Bench, we developed an intuitive Gradio-based annotation interface. Annotators
162 were shown images sampled from each category and asked to write natural language queries aligned
163 with the category theme. These queries were then evaluated using predictions from three anonymized
164 MLLMs. If one or fewer models produced a correct prediction as judged by human evaluators, the
165 query was considered sufficiently challenging and accepted for inclusion in the dataset. Following
166 this, the annotators used the same interface to annotate the target points directly on the image. A SAM
167 model was used to generate initial masks based on the selected point, and users could refine these
168 masks by editing or removing portions before submission. Finally, a separate group of annotators
169 manually verified the masks to ensure they accurately reflected the user-generated queries.

170 3.3 Point-Battle

171 As MLLMs increasingly incorporate visually grounded reasoning and pointing capabilities, static
172 benchmarks become inadequate for evaluating performance in open-ended, real-world scenarios—particularly with respect to human preferences. To address this limitation, we introduce
173 **Point-Battle**, a dynamic platform for pairwise evaluation of MLLMs’ pointing abilities based on
174 user-provided language instructions. Point-Battle adopts a head-to-head evaluation format inspired by
175 Chatbot Arena [8], implemented via a Gradio-based web interface. In each round, two anonymized
176 models are randomly sampled from the top performers in **Point-Bench**—including GPT-4o, Gemini
177 2.5 Flash, Molmo-7B-D, Qwen2.5-VL-7B, and Grok-2 Vision. Users submit a natural language
178 instruction and select an image from a curated dataset (post-April 20, 2025) or upload their own. The
179 two models return point predictions, which are displayed side by side. Participants vote for the better
180 output or select both good” or both bad” if applicable. No preset prompts are provided, encouraging
181 diverse and unbiased instructions. Model identities are kept anonymous to prevent bias. Since launch,
182 Point-Battle has collected over 4,500 votes from approximately 100 participants worldwide. Unlike
183 the static **Point-Bench**, which may be subject to overfitting during model development, Point-Battle
184 serves as a continuously updated benchmark that captures real-time human preferences and tracks
185 progress in visually grounded reasoning across MLLMs. As Point-Battle scales, it also serves as a
186 platform for collecting pointing data.

188 3.4 Point-Act

189 The first two stages of Point Arena evaluate MLLMs’ pointing capabilities through quantitative
190 metrics and human preference assessments. However, pointing is only meaningful insofar as it
191 enables real-world utility. To evaluate such support, we introduce *Point-Act*—an interactive system
192 where users issue natural-language instructions via a GUI to a double-blind MLLM. The model
193 generates one or more predicted points, which are translated into actionable commands for an xArm
194 6 Lite robot. The robot executes a pick-or-place action at the indicated location using depth sensing
195 for spatial reasoning. This setup operationalizes pointing into end-to-end physical manipulation,
196 bridging language grounding with robotic control. Point-Act highlights the downstream consequences
197 of grounding precision: even small localization errors can cause execution failures, whereas accurate
198 predictions enable consistent real-world success.

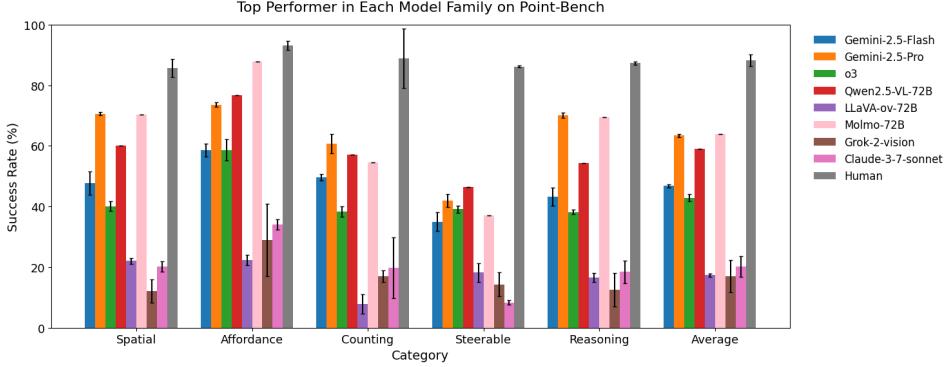


Figure 3: **Success rates of MLLMs on Point-Bench** across six task categories: *Spatial*, *Affordance*, *Counting*, *Steerable*, *Reasoning*, and *Average*. Each bar represents the mean success rate (%) for a given model, with error bars indicating standard deviation across three evaluation runs. The “Human” bar serves as an upper-bound reference. The results demonstrate substantial performance disparities, with top models (e.g., GPT-4o, Gemini-2.5-Pro, Molmo-72B) achieving near-human accuracy in select categories, while others (e.g., LLaVA, Grok, and Claude) consistently underperform.

199 4 Experiments

200 We evaluate a range of multimodal large language models (MLLMs)—both proprietary and open-
 201 source—using three components: **Point-Bench** (static benchmark evaluation), **Point-Battle** (human
 202 preference comparison), and **Point-Act** (real-world robotic execution). Section 4.1 describes the
 203 evaluation protocols, including model selection, prompting, and success metrics. Section 4.2 presents
 204 results on model performance and the impact of pointing supervision. Section 4.3 presents results
 205 demonstrating the correlation between benchmark accuracy, human preference judgments, and real-
 206 world task performance in pointing tasks. Section 4.4 includes ablations on prompt structure and
 207 output formats using GPT-4o to analyze factors affecting pointing accuracy.

208 4.1 Evaluation Setup

209 All evaluations were performed under zero-shot prompting conditions. To ensure consistent outputs
 210 across models with differing internal coordinate systems—particularly proprietary ones—we adopted
 211 a standardized output format: $[x, y]$, where x and y denote horizontal and vertical pixel coordinates,
 212 respectively. This format was used across all models, except for those like Molmo, Qwen2.5-VL, and
 213 Gemini, which provide explicit coordinate outputs or prompting instructions.

214 Success was measured using a binary metric: a prediction was considered correct if the point lay
 215 within the target mask. For non-counting tasks, models were prompted to predict a single point; if
 216 multiple were returned, only the first point was evaluated, assuming it reflected the highest-confidence
 217 prediction due to the autoregressive generation process.

218 **Point-Bench.** We benchmarked 16 MLLMs (spanning open-source and proprietary models, including
 219 key variants). Each model was evaluated on the same 982 image-instruction pairs, three times
 220 independently, to compute means and standard deviations. Open-source models were executed locally
 221 on NVIDIA A100 GPUs, while proprietary models were accessed via public APIs.

222 **Point-Battle.** To measure alignment with human preferences, we released a live evaluation platform
 223 and promoted it via social media and mailing lists. Users voted on head-to-head comparisons
 224 between anonymous model outputs. Elo ratings were computed from pairwise comparisons excluding
 225 ambiguous votes (“both good” or “both bad”).

226 **Point-Act.** We recruited 10 remote participants to interact with our real-world robot setup. For a fixed
 227 scene, participants evaluated three agents—Molmo-7B-D, GPT-4o, and a human reference—across
 228 three trials. After each condition, they completed a System Usability Scale (SUS) survey.

229 **Models.** We evaluate variants from Molmo [10], Gemini [28], OpenAI [1], Claude [3], Grok [32],
 230 LLaVA [20], and Qwen [4]. See appendix for details.

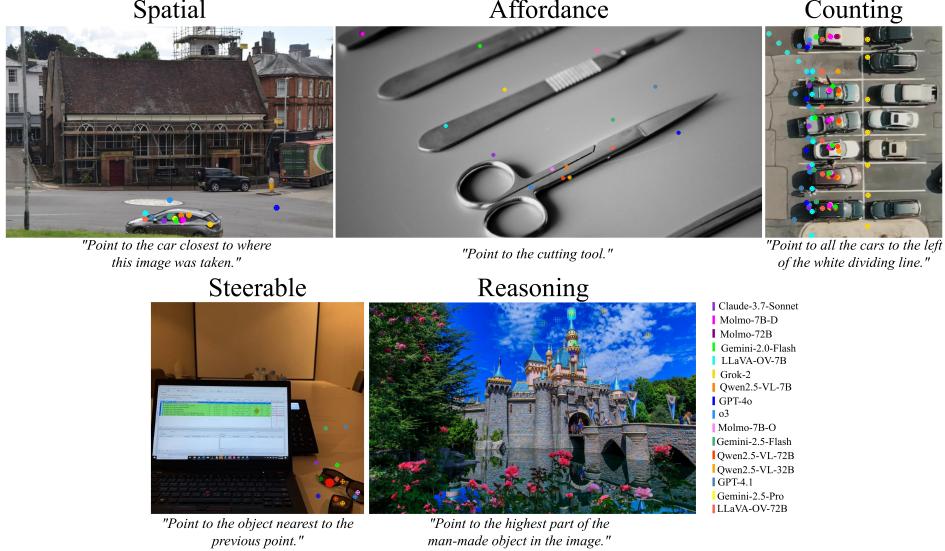


Figure 4: **Qualitative predictions across Point-Bench categories.** Example model predictions are shown for each of the five Point-Bench categories: **Spatial**, **Affordance**, **Counting**, **Steerable**, and **Reasoning**. Each colored dot corresponds to a prediction from a different MLLM, labeled by model name in the legend. These examples highlight the diversity of pointing behaviors and the variation in performance across models.

231 4.2 Main Results

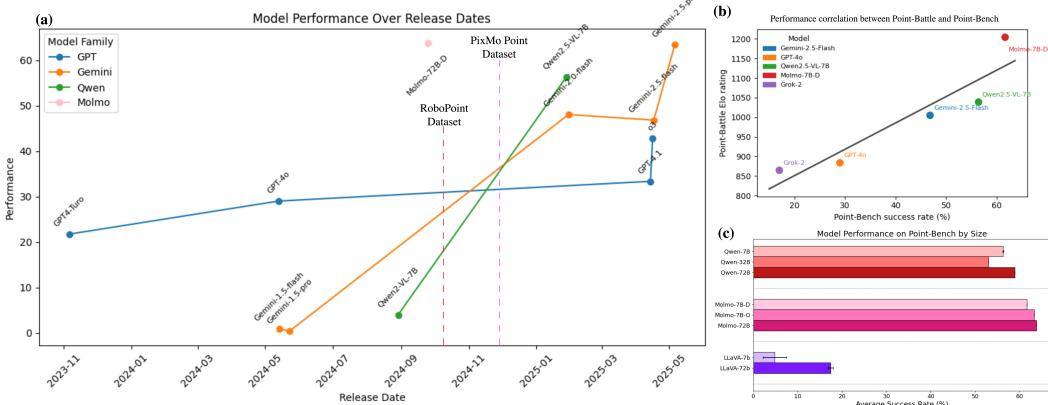


Figure 5: **Insights from Point-Battle and Point-Bench.** (a) Point-Bench accuracy over time by model family. Performance jumps notably post-PixMo (dashed line, Dec 2024), with GPT-4.1 and Gemini-2.0-Flash showing large gains over predecessors, hinting at pointing supervision. (b) Point-Battle and Point-Bench are strongly correlated ($R^2 = 0.85$), validating consistency. (c) Open-source model performance vs. model size shows only minor gains, indicating diminishing returns with scale.

232 **Open-source models perform comparably to proprietary models in pointing accuracy.** Point-
233 Bench results show that open-source MLLMs explicitly trained on pointing data often match or
234 outperform proprietary models. For example, Molmo-7B outperformed Gemini-2.5-Pro by 0.43
235 percentage points—a statistically insignificant margin ($p \approx 0.29$). In affordance reasoning, open-
236 source models like Molmo-7B and Qwen2.5-VL consistently exceed proprietary baselines. Overall,
237 Molmo-7B achieves the highest performance on the Point-Bench benchmark as shown in Table 4.

238 **Pointing supervision significantly boosts performance.** Access to explicit pointing data is a key
 239 driver of model accuracy as shown in Figure 5a. Within the Qwen family, incorporating the PixMo
 240 corpus into Qwen2.5-VL-7B increased performance to 52.3%, a substantial gain over the 17.4%
 241 achieved by Qwen2-VL-7B, which did not use such data. In contrast, LLaVA variants—also trained
 242 without explicit pointing supervision—achieved only 4.8–17.4% on average.

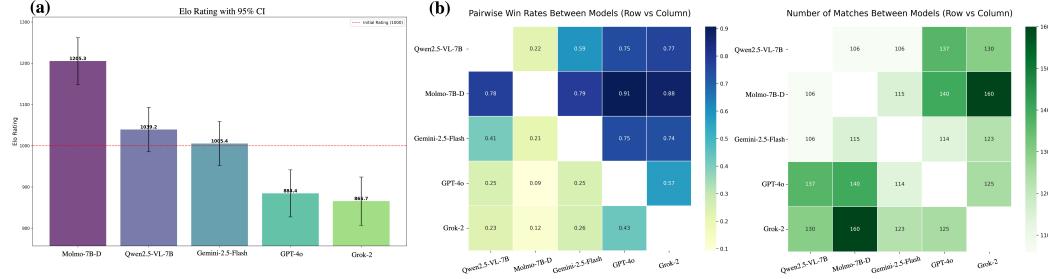


Figure 6: **Performance on Human Preference Evaluation with Point-Battle.** We collected over 4,500 votes from more than 100 global participants. Based on the Elo ratings derived from these votes, we observed a clear preference for outputs from open-source models such as Molmo-7B-D and Qwen2.5-VL-7B, which consistently outperformed proprietary models in terms of human preference.

243 **Proprietary models likely benefit from open-source pointing datasets.** While proprietary training
 244 data is opaque, we observe large performance jumps in models released shortly after the PixMo
 245 [11] and RoboPoint dataset [36]. For instance, GPT-o3 improved by 21.1 percentage points over
 246 GPT-4-Turbo, and Gemini-2.5-Flash improved by 45.9 points over Gemini-1.5-Flash (Fig-
 247 ure 5a). These results suggest that recent proprietary models may have incorporated PixMo or a
 248 similar corpus.

249 **Open-source models align more closely with human preferences.** In Point-Battle, Molmo-7B-D
 250 outperformed Gemini-2.5-Flash by 196 Elo points. Their 95% confidence intervals do not overlap,
 251 and Molmo-7B-D won 79% of the 115 direct head-to-head comparisons, as shown in Figure 6. Both
 252 Qwen2.5-VL-7B and Molmo-7B-D surpass proprietary models in human preference evaluations and
 253 exceed the 1000-point baseline, indicating a statistically significant advantage over random guessing.
 254 However, in terms of preference-aligned pointing performance, Molmo-7B-D remains clearly superior
 255 to Qwen2.5-VL-7B.

256 **Molmo excels on Point-Act evaluation.** User study results shown in Figure 7 that Molmo-7B-D out-
 257 performs the proprietary GPT-4o model by a substantial margin, achieving 65% higher performance
 258 and approaching human (oracle) baseline levels. This superiority is also reflected in user preference,
 259 with Molmo-7B-D scoring 60.3 points higher in SUS than GPT-4o.

260 **Model size does not impact pointing performance.** As shown in Figure 5c, the performance of
 261 open-source models (LLaVA-OV, Molmo, and Qwen-VL) on Point-Bench remains largely unchanged
 262 with increased model size. For example, Qwen2.5-VL-7B performs within 3% of Qwen2.5-VL-72B,
 263 and Molmo-7B-0 differs by less than 1% from Molmo-72B. These results suggest that scaling model
 264 size does not significantly improve pointing accuracy.

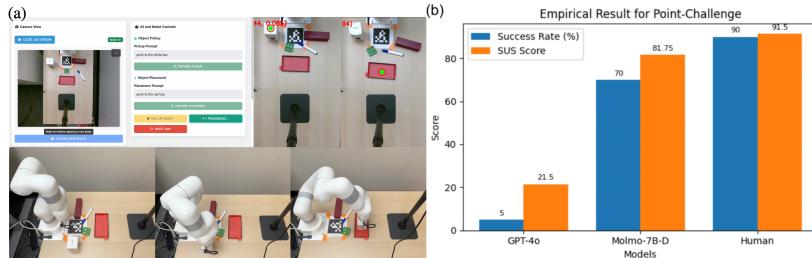


Figure 7: **Overview of the Point-Act system.** (a) The Point-Act manipulation setup enables remote control of a real-world xArm 6 Lite robot via language instructions, allowing users to evaluate pointing MLLMs. (b) User-blind evaluations and SUS preference scores collected for each model.

265 **4.3 Results between three evaluation frameworks.**

266 The three-stage evaluation of MLLMs’ pointing capabilities should not be viewed as isolated components,
267 but as complementary steps in a progressive pipeline. As MLLMs improve, they are expected
268 to advance through these stages. Therefore, understanding the correlation and agreement between
269 stages is crucial for assessing consistent performance gains.

270 **Human-preference and static dataset evaluations are highly consistent.** Point-Bench’s static
271 dataset will inevitably plateau as MLLMs improve by training on ever-larger, real or synthetic
272 pointing corpora (e.g., RoboPoint [36]). To stay ahead, we introduce **Point-Battle**, a live arena that
273 updates continuously and enables open-ended model comparison in real time. Validating this setup,
274 we re-evaluated the models tested on Point-Bench and observed strong alignment: Point-Battle scores
275 correlate with Point-Bench results at $R^2 = 0.85$ (Figure 5b).

276 **Point-Bench accuracy predicts real-world task success.** We validated Point-Bench as a reliable
277 proxy by testing three agents—Molmo-7B-D, GPT-4o, and a human reference—on Point-Act. Success
278 rates closely aligned with Point-Bench scores, yielding a strong linear correlation ($R^2 = 0.92$). This
279 high correlation indicates that Point-Bench is a reliable proxy for the pointing capability of multimodal
280 LLMs in practical settings.

281 **4.4 What Other Factors Drive Pointing Performance?**

282 To understand the design choices that impact pointing, we conducted ablations on GPT-4o using
283 variations in prompt structure and output representation, as shown in Table 1.

284 **Targeted prompts outperform verbose reasoning.** Incorporating Chain-of-Thought (CoT) reasoning
285 reduced pointing accuracy by 2.9% for GPT-4o and by a substantial 16% for Gemini-2.5-Flash.
286 Using raw, unfiltered user queries led to an additional drop of 2.6% and 3.7% for GPT-4o and
287 Gemini-2.5-Flash, respectively. These results suggest that clear, targeted prompts with well-
288 defined coordinate systems are crucial for effective pointing, while additional reasoning through
289 language does not enhance MLLMs’ pointing capabilities.

Table 1: Performance (%) of GPT-4o and Gemini-2.5-Flash variants across evaluation categories.

Method	Affordance	Spatial	Reasoning	Steerability	Counting	Average
GPT-4o + In-Context (2-shots)	46.0	26.7	23.9	22.5	33.2	30.4
GPT-4o + Chain-of-Thought (CoT) [31]	41.4	24.6	21.2	13.5	32.1	26.6
GPT-4o + Unparsed language instruction	37.8	23.7	19.7	21.9	31.1	26.9
GPT-4o (Default)	42.4	25.6	23.8	24.5	31.1	29.5
Gemini-2.5-Flash + In-Context (2-shots)	50.0	40.5	35.8	32.0	40.3	39.7
Gemini-2.5-Flash + Chain-of-Thought (CoT) [31]	44.4	25.6	24.9	26.0	33.7	30.9
Gemini-2.5-Flash + Unparsed language instruction	55.0	44.6	43.5	25.0	47.5	43.1
Gemini-2.5-Flash (Default)	58.6	47.7	43.2	35.0	49.7	46.8

290 **5 Limitations, Discussion, and Conclusions**

291 **Discussions.** PointArena benchmarks spatial reasoning in multimodal models via static datasets
292 (Point-Bench) and human preference comparisons (Point-Battle). To improve annotation quality,
293 we propose replacing grid-based tools with free-form contouring for finer object boundaries. We
294 also augment Point-Bench with user-generated Point-Battle data to reduce staleness, and introduce
295 adaptive sampling to compare similarly performing models for more informative evaluations.

296 **Limitations.** Current annotations rely on SAM [19] and grid refinement, which often yield coarse
297 masks—especially for fine-grained objects—reducing quality. Static benchmarks also risk data
298 leakage as models train on public data. Additionally, Point-Battle’s uniform sampling wastes
299 comparisons between mismatched models.

300 **Conclusion.** PointArena provides a scalable framework for evaluating grounded pointing, blending
301 controlled benchmarks with live user feedback. While limitations remain in annotation quality,
302 dataset freshness, and evaluation efficiency, our proposed updates improve precision, scalability, and
303 signal strength for real-world spatial reasoning assessments

304 **References**

- 305 [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
306 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4
307 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 308 [2] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas J. Guibas.
309 ReferIt3D: Neural listeners for fine-grained 3d object identification in real-world scenes. In
310 *16th European Conference on Computer Vision (ECCV)*, 2020.
- 311 [3] Anthropic. The claude 3 model family: Opus, sonnet, haiku. <https://www.anthropic.com/news/clause-3-family>, 2024. Claude-3 Model Card.
- 312 [4] Shuai Bai, Kebin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,
313 Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*,
314 2025.
- 315 [5] Suma Bailis, Jane Friedhoff, and Feiyang Chen. Werewolf arena: A case study in llm evaluation
316 via social deduction, 2024.
- 317 [6] Jeffrey P Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C Miller,
318 Robin Miller, Aubrey Tatarowicz, Brandy White, Samuel White, et al. Vizwiz: nearly real-time
319 answers to visual questions. In *Proceedings of the 23nd annual ACM symposium on User*
320 *interface software and technology*, pages 333–342, 2010.
- 321 [7] Dave Zhenyu Chen, Angel X. Chang, and Matthias Nießner. Scanrefer: 3d object localization
322 in rgb-d scans using natural language, 2020.
- 323 [8] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolaos Angelopoulos, Tianle Li,
324 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica.
325 Chatbot arena: An open platform for evaluating llms by human preference, 2024.
- 326 [9] Harm de Vries, Florian Strub, Sarath Chandar, Olivier Pietquin, Hugo Larochelle, and Aaron
327 Courville. Guesswhat?! visual object discovery through multi-modal dialogue, 2017.
- 328 [10] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park,
329 Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson,
330 Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris
331 Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert,
332 Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat,
333 Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon
334 Borchardt, Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa
335 Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi,
336 Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and
337 open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*, 2024.
- 338 [11] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park,
339 Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson,
340 Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris
341 Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert,
342 Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat,
343 Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon
344 Borchardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick,
345 Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick,
346 Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for
347 state-of-the-art vision-language models, 2024.
- 348 [12] Jiafei Duan, Yi Ru Wang, Mohit Shridhar, Dieter Fox, and Ranjay Krishna. Ar2-d2: Training a
349 robot without a robot. *arXiv preprint arXiv:2306.13818*, 2023.
- 350 [13] Jiafei Duan, Wentao Yuan, Wilbert Pumacay, Yi Ru Wang, Kiana Ehsani, Dieter Fox, and
351 Ranjay Krishna. Manipulate-anything: Automating real-world robots using vision-language
352 models. *arXiv preprint arXiv:2406.18915*, 2024.

- 354 [14] Petko Georgiev et al. Gemini 1.5: Unlocking multimodal understanding across millions of
 355 tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- 356 [15] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck,
 357 and Aleksandra Faust. A real-world webagent with planning, long context understanding, and
 358 program synthesis. *arXiv preprint arXiv:2307.12856*, 2023.
- 359 [16] Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith,
 360 and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal
 361 language models. *arXiv preprint arXiv:2406.09403*, 2024.
- 362 [17] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring
 363 to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical
 364 methods in natural language processing (EMNLP)*, pages 787–798, 2014.
- 365 [18] Richard Kelley and Duncan Wilson. Tournament evaluation of large language models, 2025.
- 366 [19] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,
 367 Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.
 368 Segment anything, 2023.
- 369 [20] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan
 370 Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint
 371 arXiv:2408.03326*, 2024.
- 372 [21] Zirui Liu, Jiatong Li, Yan Zhuang, Qi Liu, Shuanghong Shen, Jie Ouyang, Mingyue Cheng, and
 373 Shijin Wang. am-elo: A stable framework for arena-based llm evaluation, 2025.
- 374 [22] Rui Min, Tianyu Pang, Chao Du, Qian Liu, Minhao Cheng, and Min Lin. Improving your
 375 model ranking on chatbot arena by vote rigging, 2025.
- 376 [23] Taiki Miyanishi, Fumiya Kitamori, Shuhei Kurita, Jungdae Lee, Motoaki Kawanabe, and
 377 Nakamasa Inoue. Cityrefer: Geography-aware 3d visual grounding dataset on city-scale point
 378 cloud data, 2023.
- 379 [24] Soroush Nasiriany, Fei Xia, Wenhao Yu, Ted Xiao, Jacky Liang, Ishita Dasgupta, Annie Xie,
 380 Danny Driess, Ayzaan Wahid, Zhuo Xu, Quan Vuong, Tingnan Zhang, Tsang-Wei Edward Lee,
 381 Kuang-Huei Lee, Peng Xu, Sean Kirmani, Yuke Zhu, Andy Zeng, Karol Hausman, Nicolas
 382 Heess, Chelsea Finn, Sergey Levine, and Brian Ichter. Pivot: Iterative visual prompting elicits
 383 actionable knowledge for vlms, 2024.
- 384 [25] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor
 385 Darrell, and Marcus Rohrbach. Multimodal explanations: Justifying decisions and pointing to
 386 the evidence, 2018.
- 387 [26] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and
 388 Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer
 389 image-to-sentence models, 2016.
- 390 [27] SYV-AI. Openarena: An open platform for llm-as-a-judge evaluation. <https://github.com/syv-ai/OpenArena>, 2024. Accessed: 2025-05-09.
- 391 [28] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat
 392 Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza,
 393 Michiel Blokzijl, Steven Bohez, Konstantinos Bousmalis, Anthony Brohan, Thomas Buschmann,
 394 Arunkumar Byravan, Serkan Cabi, Ken Caluwaerts, Federico Casarini, Oscar Chang, Jose Enrique
 395 Chen, Xi Chen, Hao-Tien Lewis Chiang, Krzysztof Choromanski, David D'Ambrosio,
 396 Sudeep Dasari, Todor Davchev, Coline Devin, Norman Di Palo, Tianli Ding, Adil Dostmohamed,
 397 Danny Driess, Yilun Du, Debidatta Dwibedi, Michael Elabd, Claudio Fantacci, Cody Fong,
 398 Erik Frey, Chuyuan Fu, Marissa Giustina, Keerthana Gopalakrishnan, Laura Graesser,
 399 Leonard Hasenclever, Nicolas Heess, Brandon Hernaez, Alexander Herzog, R. Alex Hofer, Jan
 400 Humplik, Atil Iscen, Mithun George Jacob, Deepali Jain, Ryan Julian, Dmitry Kalashnikov,
 401 M. Emre Karagozler, Stefani Karp, Chase Kew, Jerad Kirkland, Sean Kirmani, Yuheng Kuang,

- 403 Thomas Lampe, Antoine Laurens, Isabel Leal, Alex X. Lee, Tsang-Wei Edward Lee, Jacky
 404 Liang, Yixin Lin, Sharath Maddineni, Anirudha Majumdar, Assaf Hurwitz Michaely, Robert
 405 Moreno, Michael Neunert, Francesco Nori, Carolina Parada, Emilio Parisotto, Peter Pastor,
 406 Acorn Pooley, Kanishka Rao, Krista Reymann, Dorsa Sadigh, Stefano Saliceti, Pannag Sanketi,
 407 Pierre Sermanet, Dhruv Shah, Mohit Sharma, Kathryn Shea, Charles Shu, Vikas Sindhwani,
 408 Sumeet Singh, Radu Soricu, Jost Tobias Springenberg, Rachel Sterneck, Razvan Surdulescu,
 409 Jie Tan, Jonathan Tompson, Vincent Vanhoucke, Jake Varley, Grace Vesom, Giulia Vezzani,
 410 Oriol Vinyals, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Fei Xia, Ted Xiao, Annie Xie,
 411 Jinyu Xie, Peng Xu, Sichun Xu, Ying Xu, Zhuo Xu, Yuxiang Yang, Rui Yao, Sergey Yaroshenko,
 412 Wenhao Yu, Wentao Yuan, Jingwei Zhang, Tingnan Zhang, Allan Zhou, and Yuxiang Zhou.
 413 Gemini robotics: Bringing ai into the physical world, 2025.
- 414 [29] Michael Tomasello, Malinda Carpenter, and Ulf Liszkowski. A new look at infant pointing.
 415 *Child development*, 78(3):705–722, 2007.
- 416 [30] Jiaqi Wang, Hanqi Jiang, Yiheng Liu, Chong Ma, Xu Zhang, Yi Pan, Mengyuan Liu, Peiran
 417 Gu, Sichen Xia, Wenjun Li, Yutong Zhang, Zihao Wu, Zhengliang Liu, Tianyang Zhong, Bao
 418 Ge, Tuo Zhang, Ning Qiang, Xintao Hu, Xi Jiang, Xin Zhang, Wei Zhang, Dinggang Shen,
 419 Tianming Liu, and Shu Zhang. A comprehensive review of multimodal large language models:
 420 Performance and challenges across different tasks, 2024.
- 421 [31] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
 422 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
 423 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 424 [32] xAI. Grok-2 model card. <https://x.ai/news/grok-2>, 2024. Large language model.
- 425 [33] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey
 426 on multimodal large language models. *National Science Review*, 11(12):nwa403, 11 2024.
- 427 [34] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. Modeling
 428 context in referring expressions. *arXiv preprint arXiv:1608.00272*, 2016.
- 429 [35] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. Modeling
 430 context in referring expressions, 2016.
- 431 [36] Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan
 432 Murali, Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial
 433 affordance prediction for robotics. *arXiv preprint arXiv:2406.10721*, 2024.
- 434 [37] Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan
 435 Murali, Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial
 436 affordance prediction in robotics. In *Proc. Conference on Robot Learning (CoRL)*, volume 270,
 437 pages 4005–4020, 2025.
- 438 [38] Ruochen Zhao, Wenxuan Zhang, Yew Ken Chia, Weiwen Xu, Deli Zhao, and Lidong Bing.
 439 Auto-arena: Automating llm evaluations with agent peer battles and committee discussions,
 440 2024.
- 441 [39] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
 442 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
 443 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.

444 **NeurIPS Paper Checklist**

445 **1. Claims**

446 Question: Do the main claims made in the abstract and introduction accurately reflect the
447 paper's contributions and scope?

448 Answer: [Yes]

449 Justification:[NA]

450 Guidelines:

- 451 • The answer NA means that the abstract and introduction do not include the claims
452 made in the paper.
- 453 • The abstract and/or introduction should clearly state the claims made, including the
454 contributions made in the paper and important assumptions and limitations. A No or
455 NA answer to this question will not be perceived well by the reviewers.
- 456 • The claims made should match theoretical and experimental results, and reflect how
457 much the results can be expected to generalize to other settings.
- 458 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
459 are not attained by the paper.

460 **2. Limitations**

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

462 Answer: [Yes]

463 Justification: [NA]

464 Guidelines:

- 465 • The answer NA means that the paper has no limitation while the answer No means that
466 the paper has limitations, but those are not discussed in the paper.
- 467 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 468 • The paper should point out any strong assumptions and how robust the results are to
469 violations of these assumptions (e.g., independence assumptions, noiseless settings,
470 model well-specification, asymptotic approximations only holding locally). The authors
471 should reflect on how these assumptions might be violated in practice and what the
472 implications would be.
- 473 • The authors should reflect on the scope of the claims made, e.g., if the approach was
474 only tested on a few datasets or with a few runs. In general, empirical results often
475 depend on implicit assumptions, which should be articulated.
- 476 • The authors should reflect on the factors that influence the performance of the approach.
477 For example, a facial recognition algorithm may perform poorly when image resolution
478 is low or images are taken in low lighting. Or a speech-to-text system might not be
479 used reliably to provide closed captions for online lectures because it fails to handle
480 technical jargon.
- 481 • The authors should discuss the computational efficiency of the proposed algorithms
482 and how they scale with dataset size.
- 483 • If applicable, the authors should discuss possible limitations of their approach to
484 address problems of privacy and fairness.
- 485 • While the authors might fear that complete honesty about limitations might be used by
486 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
487 limitations that aren't acknowledged in the paper. The authors should use their best
488 judgment and recognize that individual actions in favor of transparency play an impor-
489 tant role in developing norms that preserve the integrity of the community. Reviewers
490 will be specifically instructed to not penalize honesty concerning limitations.

491 **3. Theory assumptions and proofs**

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

494 Answer: [Yes]

495 Justification: [NA]

496 Guidelines:

- 497 • The answer NA means that the paper does not include theoretical results.
- 498 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 499 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 500 • The proofs can either appear in the main paper or the supplemental material, but if 501 they appear in the supplemental material, the authors are encouraged to provide a short 502 proof sketch to provide intuition.
- 503 • Inversely, any informal proof provided in the core of the paper should be complemented 504 by formal proofs provided in appendix or supplemental material.
- 505 • Theorems and Lemmas that the proof relies upon should be properly referenced.

506 **4. Experimental result reproducibility**

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

511 Answer:[Yes]

512 Justification: [NA]

513 Guidelines:

- 514 • The answer NA means that the paper does not include experiments.
- 515 • If the paper includes experiments, a No answer to this question will not be perceived 516 well by the reviewers: Making the paper reproducible is important, regardless of 517 whether the code and data are provided or not.
- 518 • If the contribution is a dataset and/or model, the authors should describe the steps taken 519 to make their results reproducible or verifiable.
- 520 • Depending on the contribution, reproducibility can be accomplished in various ways.
For example, if the contribution is a novel architecture, describing the architecture fully 521 might suffice, or if the contribution is a specific model and empirical evaluation, it may 522 be necessary to either make it possible for others to replicate the model with the same 523 dataset, or provide access to the model. In general, releasing code and data is often 524 one good way to accomplish this, but reproducibility can also be provided via detailed 525 instructions for how to replicate the results, access to a hosted model (e.g., in the case 526 of a large language model), releasing of a model checkpoint, or other means that are 527 appropriate to the research performed.
- 528 • While NeurIPS does not require releasing code, the conference does require all submissions 529 to provide some reasonable avenue for reproducibility, which may depend on the 530 nature of the contribution. For example
 - 531 (a) If the contribution is primarily a new algorithm, the paper should make it clear how 532 to reproduce that algorithm.
 - 533 (b) If the contribution is primarily a new model architecture, the paper should describe 534 the architecture clearly and fully.
 - 535 (c) If the contribution is a new model (e.g., a large language model), then there should 536 either be a way to access this model for reproducing the results or a way to reproduce 537 the model (e.g., with an open-source dataset or instructions for how to construct 538 the dataset).
 - 539 (d) We recognize that reproducibility may be tricky in some cases, in which case 540 authors are welcome to describe the particular way they provide for reproducibility.
In the case of closed-source models, it may be that access to the model is limited in 541 some way (e.g., to registered users), but it should be possible for other researchers 542 to have some path to reproducing or verifying the results.

545 **5. Open access to data and code**

546 Question: Does the paper provide open access to the data and code, with sufficient instruc-
547 tions to faithfully reproduce the main experimental results, as described in supplemental
548 material?

549 Answer: [Yes]

550 Justification: [NA]

551 Guidelines:

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

571 **6. Experimental setting/details**

572 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-

573 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the

574 results?

575 Answer: [Yes]

576 Justification: [NA]

577 Guidelines:

- 578 • The answer NA means that the paper does not include experiments.
- 579 • The experimental setting should be presented in the core of the paper to a level of detail
- 580 that is necessary to appreciate the results and make sense of them.
- 581 • The full details can be provided either with the code, in appendix, or as supplemental
- 582 material.

583 **7. Experiment statistical significance**

584 Question: Does the paper report error bars suitably and correctly defined or other appropriate

585 information about the statistical significance of the experiments?

586 Answer: [Yes]

587 Justification: [NA]

588 Guidelines:

- 589 • The answer NA means that the paper does not include experiments.
- 590 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
- 591 dence intervals, or statistical significance tests, at least for the experiments that support
- 592 the main claims of the paper.
- 593 • The factors of variability that the error bars are capturing should be clearly stated (for
- 594 example, train/test split, initialization, random drawing of some parameter, or overall
- 595 run with given experimental conditions).
- 596 • The method for calculating the error bars should be explained (closed form formula,
- 597 call to a library function, bootstrap, etc.)
- 598 • The assumptions made should be given (e.g., Normally distributed errors).
- 599 • It should be clear whether the error bar is the standard deviation or the standard error
- 600 of the mean.

- 601 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 602 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 603 of Normality of errors is not verified.
 604 • For asymmetric distributions, the authors should be careful not to show in tables or
 605 figures symmetric error bars that would yield results that are out of range (e.g. negative
 606 error rates).
 607 • If error bars are reported in tables or plots, The authors should explain in the text how
 608 they were calculated and reference the corresponding figures or tables in the text.

609 8. Experiments compute resources

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

613 Answer: [Yes]

614 Justification: [NA]

615 Guidelines:

- 616 • The answer NA means that the paper does not include experiments.
- 617 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 618 or cloud provider, including relevant memory and storage.
- 619 • The paper should provide the amount of compute required for each of the individual
 620 experimental runs as well as estimate the total compute.
- 621 • The paper should disclose whether the full research project required more compute
 622 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 623 didn't make it into the paper).

624 9. Code of ethics

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

627 Answer: [Yes]

628 Justification:[NA]

629 Guidelines:

- 630 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 631 • If the authors answer No, they should explain the special circumstances that require a
 632 deviation from the Code of Ethics.
- 633 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 634 eration due to laws or regulations in their jurisdiction).

635 10. Broader impacts

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

638 Answer: [Yes]

639 Justification: [NA]

640 Guidelines:

- 641 • The answer NA means that there is no societal impact of the work performed.
- 642 • If the authors answer NA or No, they should explain why their work has no societal
 643 impact or why the paper does not address societal impact.
- 644 • Examples of negative societal impacts include potential malicious or unintended uses
 645 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 646 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 647 groups), privacy considerations, and security considerations.
- 648 • The conference expects that many papers will be foundational research and not tied
 649 to particular applications, let alone deployments. However, if there is a direct path to
 650 any negative applications, the authors should point it out. For example, it is legitimate
 651 to point out that an improvement in the quality of generative models could be used to

652 generate deepfakes for disinformation. On the other hand, it is not needed to point out
653 that a generic algorithm for optimizing neural networks could enable people to train
654 models that generate Deepfakes faster.

- 655 • The authors should consider possible harms that could arise when the technology is
656 being used as intended and functioning correctly, harms that could arise when the
657 technology is being used as intended but gives incorrect results, and harms following
658 from (intentional or unintentional) misuse of the technology.
- 659 • If there are negative societal impacts, the authors could also discuss possible mitigation
660 strategies (e.g., gated release of models, providing defenses in addition to attacks,
661 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
662 feedback over time, improving the efficiency and accessibility of ML).

663 11. Safeguards

664 Question: Does the paper describe safeguards that have been put in place for responsible
665 release of data or models that have a high risk for misuse (e.g., pretrained language models,
666 image generators, or scraped datasets)?

667 Answer: [Yes]

668 Justification: [NA]

669 Guidelines:

- 670 • The answer NA means that the paper poses no such risks.
- 671 • Released models that have a high risk for misuse or dual-use should be released with
672 necessary safeguards to allow for controlled use of the model, for example by requiring
673 that users adhere to usage guidelines or restrictions to access the model or implementing
674 safety filters.
- 675 • Datasets that have been scraped from the Internet could pose safety risks. The authors
676 should describe how they avoided releasing unsafe images.
- 677 • We recognize that providing effective safeguards is challenging, and many papers do
678 not require this, but we encourage authors to take this into account and make a best
679 faith effort.

680 12. Licenses for existing assets

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

684 Answer: [Yes]

685 Justification: [NA]

686 Guidelines:

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

702 13. New assets

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

705 Answer: [Yes]

706 Justification: [NA]

707 Guidelines:

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

716 **14. Crowdsourcing and research with human subjects**

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

720 Answer: [NA]

721 Justification: [NA]

722 Guidelines:

- 723 • The answer NA means that the paper does not involve crowdsourcing nor research with
724 human subjects.
- 725 • Including this information in the supplemental material is fine, but if the main contribu-
726 tion of the paper involves human subjects, then as much detail as possible should be
727 included in the main paper.
- 728 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
729 or other labor should be paid at least the minimum wage in the country of the data
730 collector.

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

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

737 Answer: [NA]

738 Justification: [NA]

739 Guidelines:

- 740 • The answer NA means that the paper does not involve crowdsourcing nor research with
741 human subjects.
- 742 • Depending on the country in which research is conducted, IRB approval (or equivalent)
743 may be required for any human subjects research. If you obtained IRB approval, you
744 should clearly state this in the paper.
- 745 • We recognize that the procedures for this may vary significantly between institutions
746 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
747 guidelines for their institution.
- 748 • For initial submissions, do not include any information that would break anonymity (if
749 applicable), such as the institution conducting the review.

750 **16. Declaration of LLM usage**

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

755 Answer: [NA]

756 Justification: [NA]

757 Guidelines:

- 758 • The answer NA means that the core method development in this research does not
759 involve LLMs as any important, original, or non-standard components.
- 760 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
761 for what should or should not be described.