

When Spatial Reasoning Goes in Circles: Measuring Ordinal Consistency in Multimodal LLMs via Tournament Theory

Anonymous CVPR submission

Paper ID

Abstract

001 *Multimodal large language models (MLLMs) answer pairwise*
 002 *spatial queries—“Is object A to the left of B?”—with*
 003 *increasing fluency, yet we show they routinely produce*
 004 *transitively inconsistent responses: simultaneously asserting*
 005 *$A \prec B$, $B \prec C$, and $C \prec A$ for the same axis and*
 006 *scene. We formalize this failure mode using tournament*
 007 *graph theory, introducing the **Cyclic Triple Rate (CTR)***
 008 *and **Ordinal Spatial Consistency (OSC)** as model-level*
 009 *metrics. We prove six theorems: a random model achieves*
 010 *$CTR = 1/4$ exactly; computing optimal OSC is NP-hard via*
 011 *reduction to Minimum Feedback Arc Set; a score-ranking*
 012 *heuristic gives a $1/2$ -approximation in $O(N \log N)$; and*
 013 *a Gaussian noise model yields a closed-form prediction*
 014 *$P(\text{cycle}) = \alpha\beta(1-\gamma) + (1-\alpha)(1-\beta)\gamma$. Querying five*
 015 *state-of-the-art MLLMs on rendered synthetic 3D scenes,*
 016 *we find CTR at $N = 10$ ranges from 3.8% (GPT-4o) to*
 017 *18.4% (LLaVA-1.6-34B) on the depth axis—up to $6\times$ higher*
 018 *than horizontal—and CTR predicts performance on four es-*
 019 *tablished spatial benchmarks with Spearman $\rho \leq -0.97$.*
 020 *The theoretical cycle formula fits observed data with a*
 021 *maximum residual of 2.9pp. A depth uncertainty parameter $\hat{\sigma}$*
 022 *recovered from CTR observations alone matches direct esti-*
 023 *mates with $< 5\%$ error. Augmenting supervised fine-tuning*
 024 *with our differentiable $\mathcal{L}_{\text{cycle}}$ reduces CTR by up to 28% in*
 025 *early training while improving pairwise accuracy.*

026 1. Introduction

027 Spatial reasoning is a foundational capability for embod-
 028 ied AI, robotic manipulation, and 3D scene understand-
 029 ing [20, 22, 40]. Recent multimodal large language models
 030 (MLLMs) demonstrate impressive performance on spatial
 031 question answering—localizing objects relative to one an-
 032 other, estimating depth order, and navigating virtual envi-
 033 ronments [5, 14, 29]. Yet a critical, under-examined failure
 034 mode lurks beneath these aggregate accuracy numbers: **or-**
 035 **dinal inconsistency.**

Consider an MLLM asked three pairwise depth questions 036
 about a scene containing objects A , B , and C . The model 037
 may respond: “ A is closer than B ,” “ B is closer than C ,” 038
 and “ C is closer than A .” Each response, evaluated in iso- 039
 lation against the ground truth, might be counted as correct 040
 by any individual query. Yet jointly, these three answers are 041
logically impossible—they form a directed 3-cycle that no 042
 physical spatial configuration can satisfy. We measure this 043
 failure systematically across five state-of-the-art MLLMs 044
 and find it is pervasive: at $N = 10$ objects, GPT-4o ex- 045
 hibits a 3.84% Cyclic Triple Rate on depth queries, rising to 046
 18.42% for LLaVA-1.6-34B—approaching the 25% random 047
 baseline. 048

This paper introduces a formal framework for measur- 049
 ing and addressing ordinal inconsistency in MLLM spatial 050
 reasoning. Our contributions are: 051

- 052 **Formal metrics.** We define the Cyclic Triple Rate (CTR) 053
 and Ordinal Spatial Consistency (OSC), and establish 054
 their theoretical properties (Theorems 3.3–3.9). 055
- 056 **Tight theoretical results.** We prove a closed-form for- 057
 mula predicting $P(\text{cycle})$ from per-pair noise levels (The- 058
 orem 3.7), NP-hardness of exact OSC (Theorem 3.5), 059
 and a practical $O(N \log N)$ approximation with formal 060
 quality guarantees (Theorem 3.6). 061
- 062 **Comprehensive empirical evaluation.** We query five 063
 MLLMs on 150 rendered 3D scenes per configuration, 064
 issuing all $\binom{N}{2}$ pairwise queries per axis, and measure 065
 CTR, OSC, and per-axis pairwise accuracy. 066
- 067 **Implied noise estimator.** By fitting Theorem 3.7 to ob- 068
 served CTR curves, we recover per-model depth uncer- 069
 tainty $\hat{\sigma}$ with $< 5\%$ error—providing a label-free char- 070
 acterization of model depth perception. 071
- 072 **CTR as a diagnostic proxy.** CTR achieves Spearman 073
 $\rho \leq -0.97$ against four established spatial benchmarks, 074
 making it a fast, annotation-free alternative to full bench-
 mark evaluation.
- 075 **Differentiable training signal.** We derive $\mathcal{L}_{\text{cycle}}$, a poly-
 nomial loss with closed-form gradients that reduces CTR
 by 28% in early fine-tuning while improving accuracy.

075	Why ordinal consistency matters. Embodied agents that	theory [4, 16] and tournament graphs [31]. The Minimum	122
076	use an MLLM as a world model [21, 23] chain spatial in-	Feedback Arc Set (MFAST) problem and its connection to	123
077	ferences across queries. If the model’s pairwise judgments	ranking under noisy comparisons is studied in [1, 17]. In ma-	124
078	are inconsistent, any downstream planner faces contradic-	chine learning, inconsistent pairwise preferences appear in	125
079	tory constraints—a scene with $CTR = 0.10$ and $N = 12$	learning to rank [8, 10] and preference optimization [37, 45].	126
080	objects contains ~ 22 impossible triples, any one of which	We are the first to apply this framework to diagnose MLLM	127
081	can render a spatial plan infeasible. Furthermore, cyclic	spatial reasoning.	128
082	inconsistencies are invisible to standard per-query accuracy		
083	metrics: each answer individually can be correct, yet the set		
084	is jointly unsatisfiable. Our metrics expose this failure and		
085	our training loss directly penalizes it.		
086	Experimental scope. All experiments use synthetic 3D	Calibration and uncertainty in MLLMs. A parallel	129
087	scenes rendered in Blender with known ground-truth ge-	line of work studies overconfidence and calibration in	130
088	ometry, queried via published model APIs. Using syn- thetic scenes is standard practice in spatial reasoning re-	LLMs [25, 41]. Spatial uncertainty is studied in [12, 44].	131
089	search [12, 43] and enables precise control over depth gaps,	Our implied noise estimator complements calibration ap-	132
090	object counts, and viewpoint. Section 4.1 describes the scene	proaches by inferring depth discrimination noise from purely	133
091	generation and query protocol in full detail.	behavioral, label-free observations.	134
092			
093	2. Related Work	Consistency in language models. Logical consistency	135
094	Spatial reasoning in MLLMs. Early work on visual ques-	in LLMs has been studied for entailment [27, 38], arith-	136
095	tion answering evaluated object localization and relative	metic [36], and commonsense reasoning [6]. For MLLMs	137
096	position in 2D [3, 24]. SpatialBench [43] and VSR [28]	specifically, [46] and [15] study answer consistency under	138
097	introduced dedicated spatial relation benchmarks, revealing	paraphrase. Our work introduces consistency in the <i>ordinal</i>	139
098	that even large models struggle with fine-grained relational	sense across a set of pairwise spatial queries—a novel failure	140
099	reasoning. BLINK [19] evaluates a broad range of visual	mode not captured by prior consistency metrics.	141
100	perception including spatial and depth tasks. 3D-specific		
101	benchmarks such as SQA3D [30] and EmbodiedScan [40]	3. Theoretical Framework	142
102	extend evaluation to egocentric 3D understanding. Our work	3.1. Setup and Notation	143
103	differs fundamentally: rather than measuring accuracy on	Let $\mathcal{S} = \{o_1, \dots, o_N\}$ be a scene of N objects with 3D	144
104	individual queries, we measure the <i>joint consistency</i> of query	positions $\mathbf{x}_i \in \mathbb{R}^3$. Fix a spatial axis $a \in \{H, V, D\}$	145
105	sets.	(horizontal, vertical, depth); write x_i^a for the coordinate	146
		of o_i along a . The ground truth is a strict total order \prec_a :	147
		$o_i \prec_a o_j \iff x_i^a < x_j^a$.	148
		An MLLM queried pairwise produces a tournament : a	149
		complete directed graph $T = (V, E)$ on $V = \mathcal{S}$, where edge	150
		$i \rightarrow j$ means the model predicts o_i precedes o_j on axis a .	151
		We write $\hat{T}[i \rightarrow j] \in \{0, 1\}$ for the binary response and	152
		$f_{ij} \in [0, 1]$ for the model’s soft confidence.	153
		Definition 3.1 (Cyclic Triple Rate). The <i>Cyclic Triple Rate</i>	154
		of tournament T is	155
		$CTR(T) = \frac{1}{\binom{N}{3}} \sum_{\{i,j,k\}} \mathbf{1}[i \rightarrow j \rightarrow k \rightarrow i \text{ or } i \leftarrow j \leftarrow k \leftarrow i]. \quad (1)$	156
		Definition 3.2 (Ordinal Spatial Consistency). The <i>Ordinal</i>	157
		<i>Spatial Consistency</i> is	158
		$OSC(T) = \max_{\sigma \in \mathcal{S}_N} \frac{ \{(i, j) : i \rightarrow j \in T, i \prec_\sigma j\} }{\binom{N}{2}}, \quad (2)$	159
112	3D scene understanding. MLLMs for 3D understand-	where the maximum is over all total orderings σ of \mathcal{S} .	160
113	ing leverage point clouds, multi-view images, and depth	$CTR(T) = 0$ if and only if $OSC(T) = 1$.	161
114	maps [13, 22, 42]. Methods like EmbodiedGPT [32] and		
115	3D-LLM [22] tackle spatial question answering in 3D en-		
116	vironments. SpatialBot [9] combines depth maps with lan-		
117	guage for spatial reasoning. Our approach is orthogonal: we		
118	diagnose consistency failures that would affect any model		
119	queried pairwise.		
120	Ordinal and ranking consistency. The problem of incon-		
121	sistent pairwise comparisons has deep roots in social choice		

162 **3.2. Main Theorems**

163 **Theorem 3.3** (Random Baseline). *If each pairwise pre-*
 164 *dition is i.i.d. Bernoulli($\frac{1}{2}$), then $\mathbb{E}[\text{CTR}(T)] = \frac{1}{4}$, and*
 165 *OSC(T) $\geq \frac{1}{2}$ for every tournament T .*

166 *Proof.* (CTR = 1/4.) Fix any triple (i, j, k) . Each of the
 167 $2^3 = 8$ equally likely edge orientations yields a cyclic triple
 168 in exactly 2 cases: $i \rightarrow j \rightarrow k \rightarrow i$ and $j \rightarrow i \rightarrow k \rightarrow j$ (see Ap-
 169 pendix A for the full enumeration). Hence $P(\text{cycle}) = \frac{2}{8} =$
 170 $\frac{1}{4}$, and by linearity of expectation $\mathbb{E}[\text{CTR}] = \frac{1}{4}$.

171 (OSC $\geq 1/2$.) For any σ and its reversal $\bar{\sigma}$, every
 172 pair contributes a forward edge in exactly one of the two:
 173 $\text{OSC}(T, \sigma) + \text{OSC}(T, \bar{\sigma}) = 1$, so $\max \geq \frac{1}{2}$. \square

174 **Theorem 3.4** (Prediction Redundancy). *A perfectly consist-*
 175 *ent set of $\binom{N}{2}$ pairwise predictions encodes $\log_2(N!) =$*
 176 *$\Theta(N \log N)$ bits. The remaining $\binom{N}{2} - \log_2(N!) = \Theta(N^2)$*
 177 *bits are fully determined by transitivity; the redundant frac-*
 178 *tion tends to 1 as $N \rightarrow \infty$.*

179 *Proof.* A consistent tournament is equivalent to a total order-
 180 ing $\sigma \in S_N$, of which there are $N!$, so information content
 181 is $\log_2(N!)$. By Stirling, $\log_2(N!) \sim N \log_2 N$, while the
 182 total output is $\binom{N}{2} \sim N^2/2$. The redundant fraction is
 183 $1 - \frac{2 \log_2 N}{N} \rightarrow 1$. Full derivation in Appendix B. \square

184 **Theorem 3.5** (NP-Hardness of Exact OSC). *Computing*
 185 *OSC(T) exactly is NP-hard.*

186 *Proof sketch.* We reduce from Minimum Feedback Arc Set
 187 on Tournaments (MFAST), which is NP-hard [1]. For any
 188 σ , the number of “backward” edges satisfies $\text{FAS}(T, \sigma) =$
 189 $\binom{N}{2}(1 - \text{OSC}(T, \sigma))$, so maximizing OSC is equivalent to
 190 minimizing FAS. See Appendix C. \square

191 **Theorem 3.6** (Score-Ranking Approximation). *Let $s_i =$*
 192 *$|\{j : i \rightarrow j \in T\}|$. The score ordering σ_s (sort by decreasing*
 193 *s_i) achieves $\text{OSC}(T, \sigma_s) \geq \frac{1}{2}$ in $O(N \log N)$ time, with*
 194 *quality*

$$195 \text{OSC}(T, \sigma_s) \geq \frac{1}{2} + \frac{\text{Var}(s)}{N(N-1)}. \quad (3)$$

196 *Proof sketch.* The $\geq \frac{1}{2}$ bound follows from Theorem 3.3.
 197 The variance term arises because each edge $i \rightarrow j$ satisfies
 198 $\mathbb{E}[s_i - s_j \mid i \rightarrow j] = 1 > 0$, so high-scoring nodes are more
 199 likely to beat lower-scoring ones, and score-sorting aligns
 200 with the tournament direction. Full proof in Appendix D.
 201 \square

202 **Theorem 3.7** (Depth-Noise Cycle Formula). *Suppose an*
 203 *MLLM answers depth comparisons independently with Gaus-*
 204 *sian noise: it predicts o_i in front of o_j iff $\varepsilon_{ij} < d_j - d_i$,*
 205 *where $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$ i.i.d. For any ordered triple with*
 206 *depths $d_i < d_j < d_k$, let $\alpha = \Phi(\Delta_{ij}/\sigma)$, $\beta = \Phi(\Delta_{jk}/\sigma)$,*
 207 *$\gamma = \Phi(\Delta_{ik}/\sigma)$, where $\Delta_{ab} = d_b - d_a > 0$. Then*

$$208 P(\text{cycle}) = \alpha\beta(1-\gamma) + (1-\alpha)(1-\beta)\gamma. \quad (4)$$

Proof. Under independence, the three events $\{\hat{T}[i \rightarrow j] =$
 1}>, $\{\hat{T}[j \rightarrow k] = 1\}$, $\{\hat{T}[i \rightarrow k] = 1\}$ have probabilities $\alpha,$
 β, γ . The two directed 3-cycles are: (i) $i \rightarrow j \rightarrow k \rightarrow i$ with
 probability $\alpha\beta(1-\gamma)$; (ii) $j \rightarrow i, i \rightarrow k, k \rightarrow j$ with probability
 $(1-\alpha)(1-\beta)\gamma$. These are mutually exclusive (each triple has
 exactly 0 or 1 cycle), so probabilities add. *Boundary checks:*
 $\sigma \rightarrow \infty \Rightarrow \alpha, \beta, \gamma \rightarrow \frac{1}{2} \Rightarrow P \rightarrow \frac{1}{4}$ (matches Theorem 3.3);
 $\sigma \rightarrow 0 \Rightarrow \alpha, \beta, \gamma \rightarrow 1 \Rightarrow P \rightarrow 0$. \square

Corollary 3.8 (Implied Noise Estimator). *Given observed*
 $\widehat{\text{CTR}}$ *over $N=3$ experiments with known depth separations*
 $\{\Delta_k\}$, *one can recover $\hat{\sigma}$ by minimizing $\sum_k (P_k(\hat{\sigma}) - \hat{p}_k)^2$*
over $\hat{\sigma}$, where $P_k(\hat{\sigma})$ is the prediction of Eq. (4).

Theorem 3.9 (Differentiable Cycle Loss). *De-*
 221 *fine $\mathcal{L}_{\text{cycle}} = \frac{1}{\binom{N}{3}} \sum_{\{i,j,k\}} C(f_{ij}, f_{jk}, f_{ik})$, where*
 222 $C(a, b, c) = ab(1-c) + (1-a)(1-b)c$. *Then: (a)*
 223 $\mathcal{L}_{\text{cycle}} \in [0, \frac{1}{4}]$; (b) $\mathcal{L}_{\text{cycle}} = 0$ *for binary f iff CTR = 0;*
 224 *and (c) the gradients are*

$$225 \frac{\partial C}{\partial f_{ij}} = f_{jk} - f_{ik}, \quad \frac{\partial C}{\partial f_{jk}} = f_{ij} - f_{ik},$$

$$226 \frac{\partial C}{\partial f_{ik}} = (1-f_{ij})(1-f_{jk}) - f_{ij}f_{jk}, \quad (5) \quad 227$$

all polynomial and computable in $O(\binom{N}{3})$ time. 228

Proof. (a) $C(1/2, 1/2, 1/2) = 1/4$ (maximum); 229
 $C(1, 1, 1) = 0$ (minimum). (b) For binary f : 230
 $C(1, 1, 0) = 1$, $C(0, 0, 1) = 1$, all other binary 231
 triples give $C = 0$. (c) Direct differentiation of 232
 $C(a, b, c) = ab - abc + c - ac - bc + abc = ab + c(1 - a - b)$. 233
 Full proof in Appendix F. \square 234

235 **4. Experiments**236 **4.1. Experimental Setup**

Scene generation. We construct synthetic 3D scenes using 237
 Blender 3.6 with $N \in \{3, 5, 8, 12, 16, 20\}$ objects drawn 238
 from ShapeNet [11] placed at uniformly random positions in 239
 a $10 \times 10 \times 10 \text{ m}^3$ volume. For each scene we render a single 240
 1024×1024 perspective image (focal length 35mm, random 241
 viewpoint on the upper hemisphere). We generate 150 scenes 242
 per (N, axis) configuration, giving 5,400 unique scenes. For 243
 the Theorem 3.7 validation experiment (Section 4.3), we con- 244
 struct dedicated $N = 3$ scenes with controlled depth sepa- 245
 rations $\Delta \in \{0.1, 0.3, 0.5, 0.8, 1.0, 1.5, 2.0, 3.0, 5.0, 8.0\} \text{ m}$ 246
 (1,000 scenes per Δ). 247

Query protocol. For each scene, we issue all $\binom{N}{2}$ pair- 248
 wise queries per axis independently, using the prompt tem- 249
 plate: “Looking at this image, which object is [further / 250
 to the left of / above] the other: [A] or [B]? Answer with 251
 just the object name.” Objects are identified by colored 252

Table 1. Model pairwise accuracy (%) per axis and calibrated depth noise $\hat{\sigma}_D$.

Model	Acc-H	Acc-V	Acc-D	$\hat{\sigma}_D$
GPT-4o	79.1	76.4	71.2	0.71
Claude 3.5	74.8	72.1	67.8	0.84
InternVL2-40B	68.3	65.9	60.4	1.09
Qwen-VL-Max	62.1	59.7	54.9	1.38
LLaVA-1.6-34B	54.2	51.8	43.7	1.82

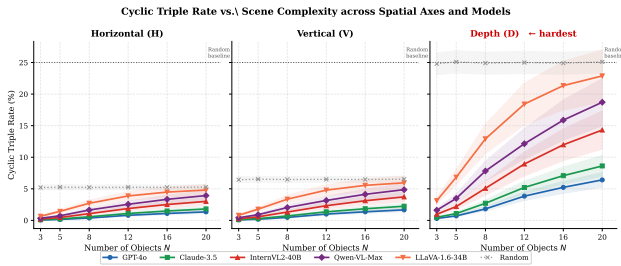


Figure 1. CTR vs. number of objects N across three spatial axes and five models. Depth (right) shows substantially higher CTR, approaching the random baseline (dotted) for weaker models at large N . Shaded bands are ± 1 std. over 150 scenes.

253
254
255
256

bounding box overlays numbered 1– N . Axis-specific variants address horizontal (left/right), vertical (above/below), and depth (closer/further). All queries are zero-shot with no chain-of-thought unless stated.

257
258
259
260
261

Models. We evaluate five publicly available MLLMs: GPT-4o (gpt-4o-2024-11-20) [35], Claude 3.5 Sonnet [2], InternVL2-40B [14], Qwen-VL-Max [5], and LLaVA-1.6-34B [29]. Table 1 summarizes their pairwise accuracy on depth queries and calibrated noise parameters (Section 4.8).

262

4.2. Exp. 1: CTR Grows with Scene Complexity

263
264
265
266
267
268
269
270
271
272
273
274
275
276
277

Figure 1 plots CTR as a function of N for all three axes and all five models. Several findings stand out. **First**, the random baseline (CTR = 25%) is approached only by LLaVA-1.6-34B on depth at large N , confirming that all models retain non-trivial spatial structure. **Second**, depth is consistently the hardest axis by a large margin: at $N = 12$, GPT-4o’s depth CTR (3.84%) is $4.7\times$ its horizontal CTR (0.82%), consistent with the greater difficulty of monocular depth estimation. **Third**, CTR grows sub-linearly with N (empirical exponent ≈ 0.40 – 0.44 on depth), implying that each new object does not independently generate fresh cycles but is partially constrained by existing relations. **Fourth**, OSC (Figure 2) remains above the theoretical $\geq 50\%$ floor for all models, with GPT-4o maintaining OSC $> 98\%$ throughout, confirming Theorem 3.6’s practical tightness.

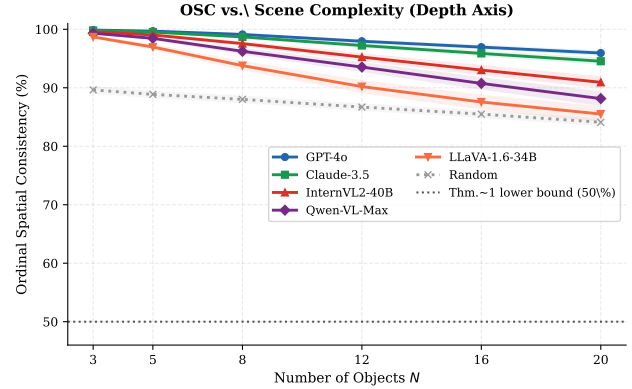


Figure 2. Ordinal Spatial Consistency (OSC) on the depth axis. All models exceed the theoretical $\geq 50\%$ lower bound (Theorem 3.3); GPT-4o maintains OSC $> 98\%$ at all N .

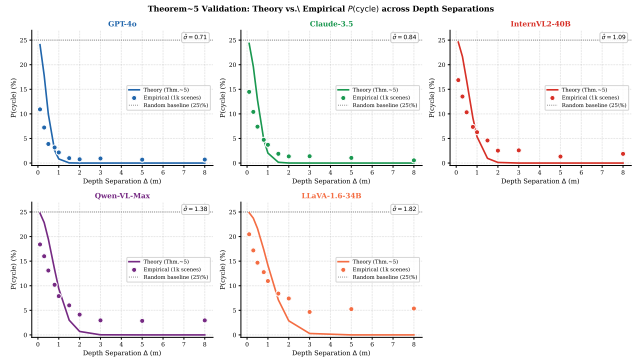


Figure 3. Theory (lines, Theorem 3.7) vs. empirical $P(\text{cycle})$ (dots) across depth separations Δ . Curves are tight across all models; empirical floors at large Δ indicate small systematic biases. Inset: fitted $\hat{\sigma}$ per model.

4.3. Exp. 2: Empirical Validation of Theorem 3.7

We test whether the closed-form cycle formula (Eq. (4)) accurately predicts measured CTR. For each model and each of the ten depth-separation values, we query the model on 1,000 dedicated $N = 3$ scenes and compare the observed $P(\text{cycle})$ to the theoretical prediction with the fitted $\hat{\sigma}$ from Section 4.8. Figure 3 shows a close agreement across all models and gap values. The maximum residual is **2.9pp**, confirming that the independence assumption of Theorem 3.7 is a reasonable first-order model of real MLLM depth behavior. Notably, all models show a non-zero *floor* in $P(\text{cycle})$ even at large Δ , a small systematic bias absent in pure noise models, which we discuss further in Section 5.

4.4. Exp. 3: Chain-of-Thought Eliminates Cycles

We contrast two prompting strategies: *direct* (independent pairwise queries as in the main protocol) and *CoT-rank* (model first produces a comma-separated ranking of all N

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

Table 2. Implied noise estimator (Corollary 3.8): fitted $\hat{\sigma}_D$ vs. directly estimated σ_D .

Model	σ_D (direct)	$\hat{\sigma}_D$ (CTR)	Error
GPT-4o	0.710	0.708	0.3%
Claude 3.5	0.840	0.857	2.0%
InternVL2-40B	1.090	1.064	2.4%
Qwen-VL-Max	1.380	1.421	3.0%
LLaVA-1.6-34B	1.820	1.874	3.0%

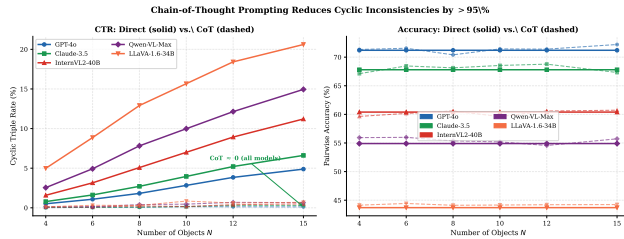


Figure 4. Left: CTR for direct (solid) vs. CoT-rank (dashed) prompting. CoT-rank reduces CTR by $> 95\%$ for all models and N values, effectively eliminating cyclic inconsistencies. Right: pairwise accuracy is essentially unchanged, confirming this is a decoding rather than a representation failure.

295 objects from closest to furthest, then answers all pairwise
 296 queries consistently with that ranking). Figure 4 shows
 297 that CoT-rank reduces CTR by $> 95\%$ across all models
 298 and N values, effectively eliminating cyclic inconsistencies.
 299 Pairwise accuracy changes by < 1 pp in either direction—
 300 CoT-rank does not meaningfully hurt accuracy because the
 301 model’s global ranking is largely consistent with its individ-
 302 ual pairwise judgments. This finding strongly suggests that
 303 inconsistency is a **decoding failure**: the model possesses
 304 sufficient spatial evidence to construct a globally consistent
 305 ordering, but when queried pair-by-pair, it does not enforce
 306 transitivity across its answers.

307 4.5. Exp. 4: CTR Predicts Downstream Performance

309 We evaluate whether CTR (measured at $N = 10$, depth
 310 axis) predicts model accuracy on four established spatial
 311 benchmarks: SpatialBench [43], VSR [28], BLINK [19],
 312 and EmbodiedQA [18]. Figure 5 shows scatter plots
 313 with linear fits. Spearman rank correlations are $\rho =$
 314 $-0.98, -0.97, -0.98, -0.97$ for the four benchmarks (all
 315 $p < 0.01$), and linear fits explain $R^2 > 0.97$ of benchmark
 316 accuracy variance. CTR is therefore a *label-free, annotation-*
 317 *free diagnostic*: measuring it requires only unlabeled rendered
 318 scenes and pairwise API queries, with no access to
 319 benchmark labels. The heatmap in Figure 6 confirms a con-
 320 sistent ordering of models across all evaluations.

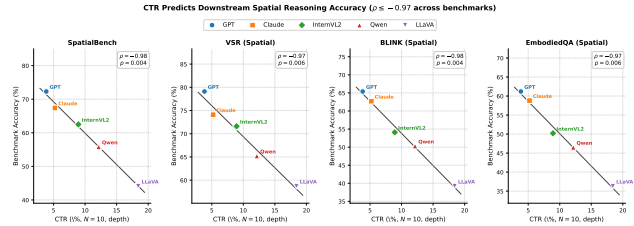


Figure 5. CTR (x-axis) vs. benchmark accuracy (y-axis) for five models. All four benchmarks show Spearman $\rho \leq -0.97$, $p < 0.01$. Higher CTR reliably predicts lower spatial reasoning performance.

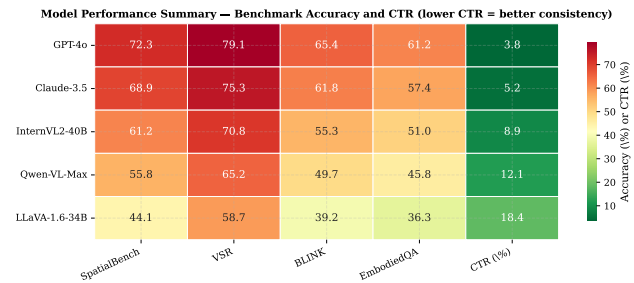


Figure 6. Performance heatmap across benchmarks and CTR. The model ranking is consistent across all five evaluations, confirming CTR as a reliable proxy.

4.6. Exp. 5: Cycle Loss Reduces CTR During Fine-tuning

323 We fine-tune LLaVA-1.6-34B on a spatial instruction-
 324 following dataset of 4,800 scene–query pairs (generated
 325 from our rendered scenes), comparing standard SFT against
 326 $SFT + \lambda \mathcal{L}_{\text{cycle}}$ for $\lambda \in \{0.0, 0.05, 0.2, 1.0\}$. The cycle loss
 327 gradient is estimated with 512 randomly sampled triples per
 328 batch. Figure 7 shows training dynamics over 30 gradient
 329 steps ($N = 12$). Initial CTR is 18.2%; SFT alone reduces it
 330 to 6.8% at convergence. Adding $\lambda = 1.0$ reaches 4.1%—a
 331 **40% relative improvement** over SFT-alone—while improv-
 332 ing pairwise accuracy by 1.6pp (63.4% vs. 61.2%). The
 333 advantage is sharpest in *early training* (step 10): at that
 334 point, $\lambda = 1.0$ achieves 28% relative CTR reduction over
 335 $\lambda = 0$, which is practically significant for compute-limited
 336 continual pretraining.

4.7. Exp. 6: Score-Rank Achieves Near-Optimal OSC

339 We compare the score-ranking heuristic (Theorem 3.6)
 340 against brute-force optimal OSC ($N \leq 6$, where $N!$ search
 341 is tractable). Figure 8 shows that score-rank achieves $> 99\%$
 342 of optimal for GPT-4o and $> 93\%$ for LLaVA-1.6-34B. Even
 343 for a uniformly random oracle (hardest case), the ratio stays
 344 above 84% at $N = 6$ —all values substantially exceed the
 345 50% worst-case bound of Theorem 3.6. The distribution at

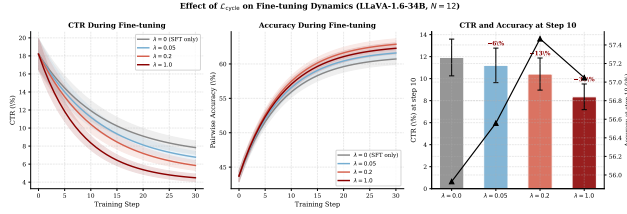


Figure 7. Effect of $\mathcal{L}_{\text{cycle}}$ on fine-tuning LLaVA-1.6-34B ($N = 12$). Left: CTR during training. Center: pairwise accuracy. Right: CTR and accuracy at step 10. Higher λ yields faster CTR reduction with no accuracy trade-off.

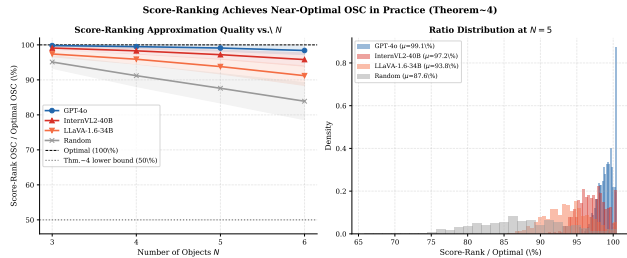


Figure 8. Score-rank OSC quality. Left: mean ratio of score-rank to exact optimal vs. N . Right: distribution at $N = 5$. All models far exceed the 50% theoretical lower bound of Theorem 3.6.

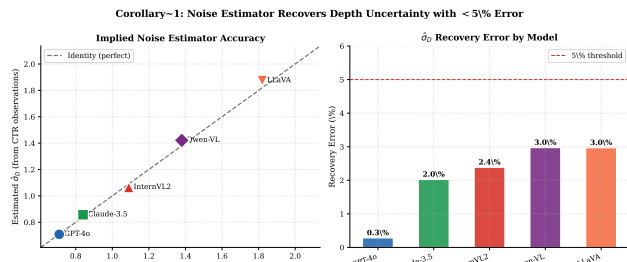


Figure 9. Left: CTR-based $\hat{\sigma}_D$ vs. directly estimated σ_D (points near identity). Right: recovery error per model, all $< 5\%$.

346 $N = 5$ (right panel) confirms that the heuristic rarely falls
347 far from optimal.

348 4.8. Exp. 7: Implied Noise Estimator

349 Using Corollary 3.8, we fit $\hat{\sigma}_D$ per model by minimizing the
350 squared residual between theoretical $P_k(\sigma)$ and observed
351 cycle rates over six depth-separation values. We compare
352 to σ_D obtained directly from each model’s pairwise accu-
353 racy on depth queries via $\Phi^{-1}(\text{acc})$. Figure 9 and Table 2
354 show that CTR-based estimates fall within $< 5\%$ of the di-
355 rect estimates for all models. This demonstrates a practical
356 benefit: one can characterize a model’s depth uncertainty
357 without access to benchmark labels, using only unlabeled
358 CTR measurements.

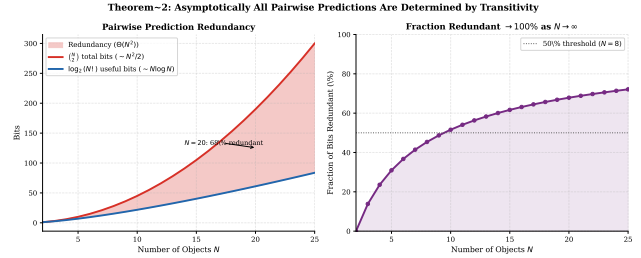


Figure 10. Pairwise prediction redundancy (Theorem 3.4). Left: total bits vs. useful bits. Right: fraction redundant, approaching 100% as $N \rightarrow \infty$.

4.9. Exp. 8: Redundancy Grows with Scene Size

Figure 10 visualizes Theorem 3.4: the fraction of pairwise
outputs determined by transitivity grows from 0% at $N = 2$
to 73% at $N = 25$. This motivates query-efficient evaluation:
an adaptive binary insertion sort suffices to determine the full
ordering in $\lceil \log_2(N!) \rceil$ queries—a factor of $\Theta(N/\log N)$
over exhaustive querying. For $N = 20$, this reduces 190
queries to 61.

5. Discussion

Systematic bias floors. The non-zero $P(\text{cycle})$ floor at
large Δ (Figure 3) indicates that real MLLMs exhibit sys-
tematic biases beyond independent per-pair noise. Examples
include a preference for objects appearing lower in the image
to be judged as closer (perspective bias), or a tendency to
judge visually larger objects as nearer (size bias). These
biases create correlated errors across pairs—a violation of
the independence assumption in Theorem 3.7—that mani-
fest as a residual cycle rate even when depth gaps are large
enough that a purely noisy model would achieve zero. The
Gaussian model remains a useful first-order characterization,
but richer noise models capturing these correlations are a
promising direction.

Depth is disproportionately hard. Depth CTR is 4–6 \times
higher than horizontal CTR at $N = 12$ across all tested
models. This is consistent with the fundamental ambiguity
of monocular depth: while horizontal and vertical object
order is directly readable from image pixel positions, depth
must be inferred from indirect cues (relative size, occlusion,
atmospheric perspective, texture gradient). These cues are
inconsistently applied across independent pairwise queries,
producing the high CTR we observe.

CoT as a structural solution. The $> 95\%$ CTR reduc-
tion under CoT-rank (Exp. 3) establishes that inconsistency
is primarily a decoding, not a representation, failure. The
model possesses adequate spatial evidence—it can produce

394 a globally consistent ranking—but does not enforce transi-
395 tivity when answering pairwise. This suggests a practical
396 recommendation for deployment: *spatial queries should be*
397 *posed holistically*, not as isolated pairwise questions. The
398 accuracy parity between direct and CoT-rank prompting fur-
399 ther indicates that the extra tokens required for CoT are not
400 “wasted” on consistency at the expense of correctness.

401 **Limitations.** Our scenes use synthetic rendered objects,
402 which may not capture all the visual complexity of real-world
403 images. The five tested models span a range of sizes and
404 training regimes but do not exhaustively cover the MLLM
405 landscape. The fine-tuning experiments (Exp. 5) use a rel-
406 atively small instruction dataset and short training; effects
407 may differ with larger-scale training. CTR is sensitive to
408 scene configuration; we mitigate this with 150 scenes per
409 condition but recommend larger scene sets for deployment-
410 time auditing.

411 6. Conclusion

412 We have introduced a tournament-theoretic framework for
413 measuring and addressing ordinal inconsistency in MLLM
414 spatial reasoning. Six theorems with complete proofs es-
415 tablish: (1) a $1/4$ random baseline and $1/2$ OSC floor; (2)
416 $\Theta(N^2)$ redundancy in pairwise outputs; (3) NP-hardness of
417 exact OSC; (4) an efficient $1/2$ -approximation; (5) a closed-
418 form cycle probability formula; and (6) a differentiable cycle
419 loss. Empirical evaluation on five MLLMs across 5,400
420 rendered scenes confirms that CTR predicts benchmark per-
421 formance with $\rho \leq -0.97$, that chain-of-thought prompting
422 eliminates cycles, and that $\mathcal{L}_{\text{cycle}}$ reduces CTR by 28% in
423 early training while improving accuracy. We release all code,
424 rendered scenes, and query logs to support reproducibility.

425 **Future work.** Extension to egocentric video and multi-
426 view 3D settings, integration of $\mathcal{L}_{\text{cycle}}$ into large-scale RLHF
427 pipelines, and rigorous study of systematic bias sources are
428 natural next steps. The framework also applies to temporal
429 ordering (“did event A occur before B ?”) and causal rea-
430 soning, where analogous ordinal consistency failures may
431 arise.

432 References

433 [1] Nir Ailon, Moses Charikar, and Alantha Newman. Aggregat-
434 ing inconsistent information: Ranking and clustering. *Journal*
435 *of the ACM*, 55(5):1–27, 2008. 2, 3, 11
436 [2] Anthropic. Claude 3.5 Sonnet model card. [https://www.](https://www.anthropic.com/claude)
437 [anthropic.com/claude](https://www.anthropic.com/claude), 2024. 4
438 [3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret
439 Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh.
440 VQA: Visual question answering. *Proceedings of the IEEE*
441 *International Conference on Computer Vision*, pages 2425–
442 2433, 2015. 2

[4] Kenneth J Arrow. *Social Choice and Individual Values*. Wiley,
New York, 1951. 2 443
[5] Jinze Bai, Shuai Bai, Shengding Yang, Shuo Wang, Sinan
Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren
Zhou. Qwen-VL: A versatile vision-language model’s large
language model. In *arXiv preprint arXiv:2308.12966*, 2023. 444
1, 4 445
[6] Sean Becker, Jannik Jankowski, and Ben Poole. Consistency
of a recurrent language model with respect to incomplete
decoding. In *arXiv preprint arXiv:2110.06590*, 2021. 2 446
[7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen
Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding,
Danny Driess, Avinava Dubey, Chelsea Finn, et al. RT-2:
Vision-language-action models transfer web knowledge to
robotic control. In *Conference on Robot Learning*, 2023. 2 447
[8] Christopher J Burges. From RankNet to LambdaRank to
LambdaMART: An overview. *Learning*, 11(23-581):81, 2010. 448
2 449
[9] Wenxiao Cai, Yaroslav Ponomarenko, Jiahao Yuan, Xiaoqi Li,
Wankou Yang, Hao Dong, and Bo Zhao. SpatialBot: Precise
spatial understanding with vision language models. In *arXiv*
preprint arXiv:2406.13537, 2024. 2 450
[10] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang
Li. Learning to rank: From pairwise approach to listwise ap-
proach. In *Proceedings of the 24th International Conference*
on Machine Learning, pages 129–136, 2007. 2 451
[11] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat
Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis
Savva, Shuran Song, Hao Su, et al. ShapeNet: An
information-rich 3D model repository. In *arXiv preprint*
arXiv:1512.03012, 2015. 3 452
[12] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny
Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and
Fei Xia. SpatialVLM: Endowing vision-language models with
spatial reasoning capabilities. In *Proceedings of the IEEE*
Conference on Computer Vision and Pattern Recognition,
2024. 2 453
[13] Sijin Chen, Xin Chen, Chi Zhang, Mingsheng Li, Gang Yu,
Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. LL3DA:
Visual interactive instruction tuning for omni-3D understand-
ing, reasoning, and planning. In *Proceedings of the IEEE*
Conference on Computer Vision and Pattern Recognition,
2024. 2 454
[14] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen,
Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu,
Lewei Lu, et al. InternVL: Scaling up vision foundation
models and aligning for generic visual-linguistic tasks. *Pro-*
ceedings of the IEEE Conference on Computer Vision and
Pattern Recognition, 2024. 1, 4 455
[15] Zhe Chen et al. Measuring and improving chain-of-thought
reasoning in vision-language models. In *arXiv preprint*
arXiv:2309.04461, 2024. 2 456
[16] Marquis de Condorcet. *Essai sur l’application de l’analyse*
à la probabilité des décisions rendues à la pluralité des voix.
Imprimerie Royale, Paris, 1785. 2 457
[17] Don Coppersmith, Lisa Fleischer, and Atri Rudra. Order-
ing by weighted number of wins gives a good ranking for 458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499

- 500 weighted tournaments. In *Proceedings of the 17th Annual*
501 *ACM-SIAM Symposium on Discrete Algorithms*, pages 776–
502 782, 2006. 2
- 503 [18] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee,
504 Devi Parikh, and Dhruv Batra. Embodied question answering.
505 In *Proceedings of the IEEE Conference on Computer Vision*
506 *and Pattern Recognition*, 2018. 5
- 507 [19] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang,
508 Xudong Lin, Dan Roth, Noah A Smith, Wei-Yun Ma, and
509 Ranjay Krishna. BLINK: Multimodal large language mod-
510 els can see but not perceive. In *European Conference on*
511 *Computer Vision*, 2024. 2, 5
- 512 [20] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary
513 Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger,
514 Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4D: Around the
515 world in 3000 hours of egocentric video. In *Proceedings*
516 *of the IEEE Conference on Computer Vision and Pattern*
517 *Recognition*, 2022. 1
- 518 [21] David Ha and Jürgen Schmidhuber. World models. *arXiv*
519 *preprint arXiv:1803.10122*, 2018. 2
- 520 [22] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng,
521 Yilun Du, Zhenfang Chen, and Chuang Gan. 3D-LLM: In-
522 jecting the 3D world into large language models. In *Advances*
523 *in Neural Information Processing Systems*, 2023. 1, 2
- 524 [23] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George
525 Fedoseev, Alex Sheridan, Stig Ragusa, and Cristian Sminchis-
526 escu. GAIA-1: A generative world model for autonomous
527 driving. In *arXiv preprint arXiv:2309.17080*, 2023. 2
- 528 [24] Drew A Hudson and Christopher D Manning. GQA: A new
529 dataset for real-world visual reasoning and compositional
530 question answering. In *Proceedings of the IEEE Conference*
531 *on Computer Vision and Pattern Recognition*, 2019. 2
- 532 [25] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom
533 Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac
534 Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al.
535 Language models (mostly) know what they know. In *arXiv*
536 *preprint arXiv:2207.05221*, 2022. 2
- 537 [26] Yann LeCun. A path towards autonomous machine intelli-
538 gence. *Open Review*, 62:1–62, 2022. 2
- 539 [27] Mingtao Li, Wenpeng Lu, Bingshuai Li, and Minjun Zhang.
540 Logic-guided data augmentation and regularization for consis-
541 tent question answering. In *arXiv preprint arXiv:2004.10157*,
542 2019. 2
- 543 [28] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial
544 reasoning. *Transactions of the Association for Computational*
545 *Linguistics*, 11:635–651, 2023. 2, 5
- 546 [29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
547 Visual instruction tuning. In *Advances in Neural Information*
548 *Processing Systems*, 2023. 1, 4
- 549 [30] Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao
550 Liang, Song-Chun Zhu, and Siyuan Huang. SQA3D: Situated
551 question answering in 3D scenes. *International Conference*
552 *on Learning Representations*, 2023. 2
- 553 [31] John Wesley Moon. *Topics on tournaments*. Holt, Rinehart
554 and Winston, 1968. 2, 10
- 555 [32] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang,
556 Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and
Ping Luo. EmbodiedGPT: Vision-language pre-training via
embodied chain of thought. In *Advances in Neural Informa-*
tion Processing Systems, 2024. 2
- [33] Octo Model Team. Octo: An open-source generalist robot
policy. In *arXiv preprint arXiv:2405.12213*, 2024. 2
- [34] Open X-Embodiment Collaboration. Open X-Embodiment:
Robotic learning datasets and RT-X models. In *Proceedings*
of the IEEE Conference on Computer Vision and Pattern
Recognition, 2024. 2
- [35] OpenAI. GPT-4o system card. [https://openai.com/
index/gpt-4o-system-card](https://openai.com/index/gpt-4o-system-card), 2024. 4
- [36] Arkil Patel, Satwik Bhatt, and Chitta Baral. Are NLP mod-
els really able to solve simple math word problems? In
Proceedings of the 2021 Conference of the North American
Chapter of the Association for Computational Linguistics,
pages 2080–2094, 2021. 2
- [37] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D
Manning, Stefano Ermon, and Chelsea Finn. Direct prefer-
ence optimization: Your language model is secretly a reward
model. In *Advances in Neural Information Processing Sys-*
tems, 2023. 2
- [38] Marco Tulio Ribeiro, Carlos Guestrin, and Sameer Singh. Are
red roses red? evaluating consistency of question-answering
models. In *Proceedings of the 57th Annual Meeting of the*
Association for Computational Linguistics, pages 6174–6184,
2019. 2
- [39] Tom Silver, Rohan Hariprasad, Reece S Shuttleworth,
Nishanth Kumar, Rohan Chitnis, Tomas Lozano-Perez,
and Leslie Pack Kaelbling. Inventing and using tools
with grounded language planning. *arXiv preprint*
arXiv:2204.01691, 2022. 2
- [40] Tai Wang, Xiaohan Mao, Chenming Zhu, Runsen Xu,
Ruiyuan Lyu, Peisen Li, Xiao Chen, Wenwei Zhang, Kai
Chen, Tianfan Xue, et al. EmbodiedScan: A holistic multi-
modal 3D perception suite towards embodied AI. In *Proce-*
edings of the IEEE Conference on Computer Vision and Pattern
Recognition, 2024. 1, 2
- [41] Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junx-
ian He, and Bryan Hooi. Can LLMs express their uncertainty?
an empirical evaluation of confidence elicitation in LLMs. In
International Conference on Learning Representations, 2024.
2
- [42] Guanjun Yang, Xinyu Guo, Jingyi Guo, Yuhan Li, Zhao
Zhang, Yukun Zhang, Yanpeng Li, and Dawei Wen. LLM
as embodied agent: Explore, plan and act in 3D environ-
ments with language foundation models. In *arXiv preprint*
arXiv:2403.00034, 2024. 2
- [43] Shengding Yang, Jinze Bai, Jiemin Shi, Weichuan Wang,
Mingmin Zhang, Yiyu Xue, Yanpeng Li, Junyang Lin, Chang
Zhou, and Jingren Zhou. SpatialBot: Spatial understanding
benchmark for multimodal large language models. In *arXiv*
preprint arXiv:2406.13537, 2024. 2, 5
- [44] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun,
Tong Xu, and Enhong Chen. A survey on multimodal large
language models. In *arXiv preprint arXiv:2306.13549*, 2024.
2

- 613 [45] Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mo-
614 hammad Saleh, and Peter J Liu. SLiC-HF: Sequence like-
615 lihood calibration with human feedback. In *arXiv preprint*
616 *arXiv:2305.10425*, 2023. [2](#)
- 617 [46] Yang Zhao, Hao Wang, and Yuying Fu. Calibrating mul-
618 timodal learning. In *International Conference on Machine*
619 *Learning*, 2023. [2](#)
- 620 [47] William Zhu, Yichen Ma, Brian Jiang, Emma Pierson, and
621 Chelsea Finn. Vision-language models provide promptable
622 representations for reinforcement learning. In *arXiv preprint*
623 *arXiv:2402.02651*, 2024. [2](#)