

MCSBench: Probing Multimodal Conceptual Structure of Multimodal LLMs

Anonymous CVPR submission

Paper ID 31

Abstract

001 *Do multimodal LLMs (MLLMs) understand concepts struc-*
 002 *tureally like humans? Inspired by cognitive principles, we*
 003 *formalize multimodal conceptual structure (MCS) as the rela-*
 004 *tional organization grounded in concept-attribute bindings,*
 005 *inter-concept relations, and hypothetical transformations*
 006 *across vision-language space. We introduce MCSBENCH, a*
 007 *diagnostic benchmark of 661 questions across seven tasks*
 008 *spanning four cognitive levels (from perceptual grounding*
 009 *to metacognitive verification), paired with 518 structural-*
 010 *integrity probes. We further propose the Structural Align-*
 011 *ment Index (SAI), an integrity-aware metric that awards*
 012 *credit only when both the answer is correct and the underly-*
 013 *ing reasoning is structurally sound. Evaluating ~80 MLLMs,*
 014 *we find that performance degrades sharply with cognitive*
 015 *depth, and SAI exposes structural brittleness that accuracy*
 016 *alone conceals. Notably, in-context golden evidence yields*
 017 *substantially smaller integrity gains than accuracy gains on*
 018 *difficult/all questions, suggesting that retrieval augmentation*
 019 *may inflate correctness without facilitating genuine struc-*
 020 *tural reasoning. Additional analyses, such as item-response*
 021 *modeling, scaling analysis, and reasoning-graph diagnosis,*
 022 *validate that MCSBENCH provides a reliable, fine-grained*
 023 *diagnostic lens into conceptual-structure failures that exist-*
 024 *ing benchmarks overlook. We will release our dataset and*
 025 *artifacts upon acceptance.*

026 1. Introduction

027 Cognitive science holds that concepts are not isolated labels
 028 but relationally organized systems: defined by discriminative
 029 attributes, structured by similarity and contrast with confus-
 030 able neighbors, and stress-tested by counterfactual transfor-
 031 mation [1, 4, 14]. A system that genuinely understands a
 032 concept should maintain this relational structure consistently
 033 across modalities, and its correctness should be auditable
 034 against evidence, not merely coincidental [11]. Whether
 035 multimodal large language models acquire such structure
 036 remains unknown because no existing benchmark jointly
 037 probes attribute binding, relational inference, counterfactual

reasoning, and reasoning integrity. 038

We formalize multimodal conceptual structure (MCS) as 039
 this relational organization and introduce MCSBENCH (in- 040
 cluding 661 base questions across 7 tasks spanning 4 cog- 041
 nitive levels, overlaid with 518 integrity-checked questions), to- 042
 gether with our proposed Structural Alignment Index (SAI), 043
 which requires correctness to be structurally justified. 044

Evaluating ~ 80 MLLMs, we find that SAI can re- 045
 order model rankings relative to accuracy, exposing struc- 046
 turally unjustified correctness that current evaluation con- 047
 ceals. Most critically, providing models with explicit evi- 048
 dence substantially improves answer accuracy but yields 049
 near-zero integrity gains on difficult items—implying that 050
 retrieval-augmented pipelines may inflate correctness with- 051
 out instilling genuine structural reasoning. These results 052
 suggest that both current evaluation practices and evidence- 053
 grounding strategies require integrity-aware redesign. 054

Our contributions: (1) MCS formalized as a cognitively 055
 grounded evaluation construct, instantiated in MCSBENCH 056
 across four cognitive levels; (2) SAI, validated via IRT 057
 calibration, separating correctness from justified correctness; 058
 (3) We provide a diagnostic analysis suite that turns the 059
 benchmark into actionable insight through systematic analy- 060
 sis across different dimensions. 061

062 2. MCSBench

063 2.1. A Unified Relational Perspective

Human does not grasp concepts as isolated atoms. Cogni- 064
 tive science argues that concepts are internally structured: 065
 what defines a concept is not a checklist of surface fea- 066
 tures/properties but a relational organization—encoding 067
 which attributes are diagnostic and discriminative for a 068
 concept, how concepts neighbor and contrast with confus- 069
 able alternatives, and how these relations are preserved 070
 or transformed under interventions and counterfactual ed- 071
 its [1, 31, 45]. We define **multimodal conceptual structure** 072
 (MCS) as this internal relational organization expressed con- 073
 sistently across vision and language: specifically, (i) *concept-* 074
attribute role knowledge—the structured binding of dis- 075
 criminative attributes to concepts grounded in visual evi- 076

077 dence [16, 43], and (ii) *inter-concept relational organiza-*
 078 *tion*—the coherent arrangement of concepts into confus-
 079 ability neighborhoods and the relational geometry that gov-
 080 erns similarity, contrast, and structural transformation across
 081 modalities [14, 40, 44]. Critically, MCS is a property of rep-
 082 resentation: a model possesses it only if its concept–attribute
 083 bindings and inter-concept relations are mutually consistent
 084 across visual and linguistic inputs, not merely correlated
 085 with label statistics. An equally important epistemic prin-
 086 ciple governs our evaluation design: answer correctness is
 087 necessary but insufficient for evidencing MCS, since a model
 088 may reach the correct option through surface shortcuts [13]
 089 without any structurally valid reasoning. We therefore re-
 090 quire correct answers to be supported by verifiable evidence
 091 and a consistent reasoning structure—the core motivation
 092 behind our integrity-aware overlay, which decouples *correct-*
 093 *ness* from *justified correctness*.

094 Guided by these principles, MCSBENCH organizes
 095 its tasks into four cognitive levels that probe MCS with
 096 increasing operational depth: L1 *Perceptual Grounding*,
 097 L2 *Relational Inference*, L3 *Hypothetical Reasoning*, and
 098 L4 *Metacognitive Verification*. Each level is implemented as
 099 one or more task variants that hold the underlying relational
 100 structure fixed while varying the operation demanded of the
 101 model (see Table 1). Additional interdisciplinary discussion
 102 is provided in Appendix.

103 2.1.1. MCSBench Task Taxonomy

104 We define seven base tasks spanning four conceptual reason-
 105 ing categories and four cognitive levels. Let q denote the
 106 query task, c a concept, I_c its object-centric image, R the
 107 relation set, A the attribute set, and $T(\cdot)$ a templated prompt.

108 *C-A Align.* predicts the correct concept–attribute relation:
 109 $\Pr(A \mid T(q), I_c)$. *CA-Discr.* picks a^+ over distractor a^- :
 110 $\Pr(A^+ \succ A^- \mid T(q), I_c)$, distinguishing the concept from
 111 its nearest confusable neighbor [44]. *CC-Align.* follows a
 112 given relation R to select the correct compositional answer:
 113 $\Pr(A^+ \succ A^- \mid T(q), R, I_c)$; success requires tracking the
 114 *relational role* of R , not category-level similarity [14]. *CC-*
 115 *Discr.* infers both the relation type and its filler without a
 116 relation label: $\Pr(R, c' \mid T(q), I_c)$.

117 *Hypo. Reas. (Interventional)* applies an explicit do-
 118 intervention [33]: $\Pr(c' \mid T(q, \text{do}(A \leftarrow A')), I_c)$. *Hypo.*
 119 *Reas. (Counterfactual)* reasons over a fully observed in-
 120 stance under a counterfactual attribute edit: $\Pr(c_{A \leftarrow A'} =$
 121 $c^{\text{cf}} \mid T(q), R, I_c)$, with $c_{A \leftarrow A'}$ the potential-outcome vari-
 122 able from Pearl’s abduction–action–prediction.

123 *MI-Discr.* selects images satisfying a textual at-
 124 tribute/relation description from a candidate set. *MI-Match.*
 125 produces a globally consistent assignment between multiple
 126 images and multiple descriptions.

127 *Reasoning-Chain Evaluation* checks whether a model’s ex-
 128 planation is structurally valid w.r.t. question and evidence,

measuring the SAI overlay gap between answer correctness
 and justified correctness.

3. Evaluation Protocol

3.1. Structural Alignment Index (SAI)

We aim to evaluate the structural alignment of conceptual
 representations in MLLMs. Traditional accuracy metrics
 are insufficient: models frequently produce correct answers
 through flawed reasoning [10], exploit spurious shortcuts
 that intensify with scale [42], and in multimodality, over 70%
 of widely-used VQA samples can be solved by language pri-
 ors alone without visual input [19]. Accuracy thus disguises
 whether a model has genuinely and faithfully learned the
 underlying conceptual structure or merely exploited statisti-
 cal regularities. To address this, we propose the **Structural**
Alignment Index (SAI), which jointly requires both an-
 swer correctness and structural reasoning fidelity. Here, we
 introduce the probabilistic and empirical form of SAI as:

$$\text{SAI} = \mathbb{E}_{q \sim \mathcal{Q}, \tau \sim \mathcal{T}_q} \left[\Pr(\hat{y} = y_q \mid q) \Pr(\hat{t} = t_{q,\tau} \mid q, \tau) \right], \quad (1)$$

$$\widehat{\text{SAI}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{K} \sum_{k=1}^K \mathbb{I}(\hat{y}_i = y_i) \mathbb{I}(\hat{t}_{ik} = t_{ik}), \quad (2)$$

where $\mathcal{T}_q = \{\tau_{q,1}, \dots, \tau_{q,K}\}$ denotes K reasoning probes
 paired with each base question q , each targeting certain struc-
 tural relation(s) that a correct solution must preserve, with a
 gold answer $t_{q,k}$. The product structure ensures that a model
 receives nonzero credit only when it answers q correctly *and*
 demonstrates faithful structural reasoning on the associated
 probes. Empirically, SAI reduces to the multiplication of
 base accuracy and mean probe accuracy per question.

3.2. Experimental Setup

Our evaluation spans 82 multimodal LLMs across three
 groups: proprietary frontier models (*e.g.*, GPT-5.2, Claude-
 Opus-4.6, Gemini-3), a broad open-source family (*e.g.*, In-
 ternVL3.5/3, Qwen-VL3/2.5, GLM, MiniCPM), and their
 respective thinking variants. All models are evaluated un-
 der two conditions. In the standard *closed-book* setting as
 default, the model receives only the task image and ques-
 tion prompt, with no access to external knowledge. In the
open-book setting, the model is additionally provided with
 in-context evidence of: (i) web-retrieved evidence via search
 tool (web), or (ii) ground-truth evidence from our curated
 evidence (evd). We enforce greedy decoding for all models
 unless specified.

4. Experiments and Analyses

We organize our analyses as: starting from the capability
 landscape under accuracy and integrity-aware measurement

Table 1. Overview of task schema design in MCSBENCH with eight concrete task variants. Cognitive levels (L0–L4) indicate increasing reasoning complexity. L1=Perceptual Grounding, L2=Relational Inference, L3=Hypothetical Reasoning, and L4=Metacognition.

Task Group	Task Variants	Cog. Level	Task Target
Concept-Attribute	Discrimination, Alignment	L1	<i>Visual-Attribute Binding</i> : Align linguistic attribute semantics with visual evidence and discriminate confusing candidates.
Concept-Concept	Discrimination, Alignment	L2	<i>Cross-Modal Relational Inference</i> : Align inter-concept (multi-hop) relations across modalities with linguistic specifications and discriminate against confusing concepts or spurious relations.
Hypothetical Reasoning	Interventional/Counterfactual	L3	<i>Causal Relation Prediction</i> : Predict how conceptual relations shift under population-level interventions or instance-level counterfactual edits across modalities.
Multi-Image Concept	Discrimination, Matching	L2	<i>Multi-Image Relational Matching</i> : Infer many-to-many correspondences between multiple concept images and linguistic descriptions specifying concept-attribute or concept-concept relations.
Reasoning-Chain	Evaluation	L4	<i>Reasoning Chain Verification</i> : Evaluate explanatory reasoning chains for evidence faithfulness, logical correctness, and soundness.

173 (§5.1), ruling out semantic confounders (§5.2), and establish-
 174 ing reliability via IRT (§5.3); then diagnosing what helps,
 175 in-context analysis (§5.4) and scaling trend (§5.5), before
 176 testing external consistency through cross-benchmark correla-
 177 tions (§5.6).

178 5. Main Evaluation

179 5.1. Evaluation Results

180 *Overview.* Table 3 reports closed-book performance across
 181 all models. The consistent performance gaps across task
 182 axes confirm that MCSBench captures distinct structural
 183 capabilities beyond standard VQA, providing a fine-grained
 184 diagnostic that general benchmarks cannot offer.

185 **Proprietary models lead with a clear gap.** GPT-5 achieves
 186 the highest average (78.37%), and the best open-source mod-
 187 els (InternVL3-78B, Qwen2.5VL-72B; 60.06%) trail the
 188 proprietary frontier by ~ 14 pp, while many models fall
 189 below 30%.

190 **A consistent cognitive gradient: L0 > L1 > L2.** Perform-
 191 ance decreases monotonically along the cognitive ladder
 192 across nearly all models (*e.g.*, InternVL3-78B: 80.46 \rightarrow
 193 66.37 \rightarrow 25.19), with the L1 \rightarrow L2 drop sharpest for non-
 194 frontier models, indicating a persistent inability in inter-
 195 concept and hypothetical relations.

196 **Multi-image processing \neq MCS.** Several models achieve
 197 competitive multi-image scores while remaining weak on
 198 L2 (*e.g.*, MiniCPM-V-4.5: MI-Discr. 75.78%, L2 21.37%),
 199 showing the bottleneck is higher-order structured conceptual
 200 inference, not multi-image processing *per se*.

201 **CoT $\not\approx$ MCS.** Thinking variants do not consistently out-
 202 perform base counterparts (GLM-4.1V-Thinking 52.04% vs.
 203 GLM-4.5V 58.09%), confirming MCS competence is not
 204 reducible to chain-of-thought prompting.

205 Table 4 reports performance on the split with Reasoning-
 206 Chain overlay:

**SAI exposes structural brittleness hidden by Base accu-
 racy.** Across all models and settings, SAI is strictly lower
 than Base with systematically varied gap magnitude. *CC-
 Align* and *Hyp. Reas.* show the steepest Base \rightarrow SAI drops
 across all tiers, indicating unfaithful reasoning or alignment
 in higher-order relations.

Proprietary–open-source gap widens under SAI. The
 top proprietary model (Gemini-3-Pro-Preview) leads the
 best open-source counterpart (InternVL3-38B) by ~ 24
 SAI pp—wider than the Base gap alone, highlighting the
 discriminative power of SAI.

**Open-book evidence narrows but does not close the SAI
 gap.** Even with ground-truth evidence, strong open-source
 models (*e.g.*, InternVL3-78B: 57.59 SAI) still trail top pro-
 prietary closed-book scores, demonstrating the ability to
 integrate evidence into faithful reasoning matters.

223 6. Mechanistic Diagnostics

224 6.1. Recognition \neq Structure

225 6.1.1. Semantic Controls Setup

226 We run two controls to test whether MCSBench diffi-
 227 culty reduces to superficial factors: sensitivity to seman-
 228 tic label cues in the prompt, and performance dependence
 229 on category recognition. We augment instructions with
 230 label/descriptor cues (semantic-label intervention) across
 231 $N=12$ models on 661 base-task items, and separately filter
 232 evaluation to each model’s top- $k\%$ recognized categories
 233 ($k \in \{10, 20, 30, 40, 50\}$). Implementation details see Ap-
 234 pendix.

235 **Semantic label cues have negligible effect.** The overall
 236 accuracy change is small (+0.87 pp, 95% CI: [0.32, 1.46]),
 237 and an equivalence test (TOST, ± 2 pp margin) confirms prac-
 238 tical negligibility ($p_{\text{TOST}}=0.0016$). Item-level prediction
 239 patterns are stable (0/12 models show significant change un-
 240 der McNemar’s test), and task-level effects are sparse—only
 241 C-A Discr. (+5.88 pp, $p=0.019$) and C-C Align. (+4.28 pp,

242 $p=0.028$) reach significance (Table 5).
 243 **Strong recognition does not rescue performance.** Recogni-
 244 tion filtering improves recognition-adjacent tasks at small
 245 k (C-A Align. +19.4 pp at $k=10\%$), confirming the filter
 246 is selective, yet relation- and reasoning-heavy tasks do not
 247 consistently benefit and can degrade (C-C Discr. -9.8 pp at
 248 $k=10\%$ – 20% , significant after BH-FDR; Fig. 4).
 249 **Implication.** MCSBench difficulty is not attributable to se-
 250 mantic cues or recognition limitations; it genuinely probes
 251 underlying multimodal structure beyond superficial recogni-
 252 tion.

253 6.2. In-Context Diagnostics

254 We ask whether explicit evidence can close the remaining
 255 gap under open-book evaluation, providing k in-context evi-
 256 dence of three types (*Fact*, *Observation*, *Mixed*) and mea-
 257 suring performance gains over a no-evidence baseline on
 258 ALL and MEDIUM-HARD subsets. We define Evidence-
 259 Utilization Efficiency $EUE(k) = \frac{\sum_m \Delta SAI_m(k)}{\sum_m \Delta Acc_m(k)}$ to quan-
 260 tify structural integrity gains relative to accuracy gains. (1)
 261 *Evidence contribution.* Mixed evidence yields the largest
 262 mean ΔAcc , yet all types produce substantially lower ΔSAI
 263 gains (Fig. 5a): evidence improves surface-level correct-
 264 ness far more than structural reasoning. (2) *Diminish-*
 265 *ing returns.* Accuracy gains follow a saturation curve
 266 ($\Delta(k) = 0.314(1 - e^{-k/4.6})$, $R^2 = 0.981$; $k_{90} \approx 11$), while
 267 ΔSAI saturates at a markedly lower ceiling (Fig. 5b). Be-
 268 yond $k \approx 8$ – 12 , additional evidence yields negligible bene-
 269 fit for correctness and even less for integrity. (3) *Integrity*
 270 *gap.* $EUE < 1$ across all k , with efficiency collapsing on
 271 MEDIUM-HARD items (Fig. 5c): evidence acts as decision
 272 support, not *structural* integrity instiller. This bottleneck
 273 motivates our structural diagnosis below.

274 6.3. Reasoning-Graph Diagnosis of Structural Rea- 275 soning

276 We extract typed reasoning graphs from model rollouts
 277 ($N=2310$; 9 models) and compute WL-based *Stability* and
 278 *Faithfulness* scores alongside evidence-anchoring edge fea-
 279 tures, diagnosing how models reason beyond answer correct-
 280 ness.

281 Three findings emerge. First, only items that are both sta-
 282 ble and faithful achieve high accuracy (mean 0.63 for Robust
 283 vs. 0.35 for Disorganized; Fig. 6), showing that structural
 284 alignment—not consistency alone, drives success. Second,
 285 grounding edges are the strongest OLS predictor of accu-
 286 racy (std. $\beta=0.099$, $p<.001$), exceeding *Stability* ($\beta=0.049$)
 287 and *Faithfulness* ($\beta=0.042$) (Table 6b), implicating explicit
 288 evidence-to-structure binding as the key bottleneck. Third,
 289 under difficult settings, *Stability* is an unreliable proxy: its
 290 correlation with accuracy flips from $\rho=+0.29$ on easy items
 291 to $\rho=-0.14$ on hard items (Fisher $z=8.59$, $p<.05$; Table 6a),
 292 revealing a stable-but-wrong regime that motivates integrity-

aware scoring over self-consistency.

293 7. Scaling Trend

294 We provide an interpretable and testable decomposition of
 295 scaling behavior on MCSBench: how much variance is ex-
 296 plained by size, how much by training recipe/lineage, and
 297 which structural relation scale the most. We regress per-task
 298 accuracy on $\ln P$ and decompose the explained variance
 299 (Fig. 7, Tab. 7).
 300

301 **Scaling exists but explains only moderate variance.** Most
 302 tasks show positive slopes (Avg. $\beta = 7.4$), yet scale alone
 303 accounts for only a moderate fraction of variance (Avg.
 304 $R_{sc}^2 \approx 0.32$), confirming that MCSBench is not reducible to
 305 a parameter-count leaderboard.

306 **Training recipe dominates residual variance.** Adding
 307 family/lineage intercepts yields the largest incremental gain
 308 (Avg. $\Delta R_{fam}^2 \approx 0.50$, $p<.05$), raising the average fit to
 309 $R_{full}^2 \approx 0.87$. Performance gaps thus manifest primarily
 310 as recipe-level offsets rather than different scaling rates;
 311 family \times scale interactions contribute negligibly (see Ap-
 312 pendix).

313 **Scaling is task- and depth-dependent.** Slopes vary substan-
 314 tially across task groups (Fig. 7): multi-image tasks scale
 315 most strongly (MI-Match $\beta = 13.5$), while hypothetical
 316 reasoning and deeper cognitive levels are persistent bottle-
 317 necks (Hypo. Reas. $\beta = 2.1$; Cog-L3 $\beta = 2.1$), with weaker
 318 family gains (Tab. 7), indicating that higher-order structural
 319 reasoning is not reliably acquired through capacity scaling
 320 alone.

321 MCSBench exhibits clear but moderate scaling, with
 322 training recipe accounting for the majority of residual vari-
 323 ation. Critically, scaling gains are non-uniform across con-
 324 ceptual relations, making MCSBench a diagnostic tool for
 325 identifying where scale helps and where targeted training
 326 are needed.

327 8. Conclusion

328 MCSBench formalizes multimodal conceptual structure as
 329 a measurable evaluation dimension and reveals, via SAI,
 330 that correct answers frequently lack structural justification, a
 331 gap unresolved by evidence provision or scale alone. These
 332 findings argue for integrity-aware evaluation and structurally
 333 grounded system design as field priorities.

References

- 334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
- [1] Lawrence W Barsalou. Perceptual symbol systems. *Behavioral and brain sciences*, 22(4):577–660, 1999. 1, 7
- [2] R Darrell Bock and Murray Aitkin. Marginal maximum likelihood estimation of item parameters: Application of an em algorithm. *Psychometrika*, 46(4):443–459, 1981. 12
- [3] David Byrne. Complexity, configurations and cases. *Theory, culture & society*, 22(5):95–111, 2005. 7, 8
- [4] Ruth MJ Byrne. Precis of the rational imagination: How people create alternatives to reality. *Behavioral and brain sciences*, 30(5-6):439–453, 2007. 1
- [5] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *Advances in Neural Information Processing Systems*, 37:27056–27087, 2024. 7
- [6] Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang. Can pre-trained vision and language models answer visual information-seeking questions? In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14948–14968, 2023. 7, 8
- [7] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2818–2829, 2023. 10
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 10, 19
- [9] Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, et al. Agent ai: Surveying the horizons of multimodal interaction. *arXiv preprint arXiv:2401.03568*, 2024. 7
- [10] Don Fallis and Peter J Lewis. Right for the wrong reasons: common bad arguments for the correct answer to the monty hall problem. *Synthese*, 207(1):53, 2026. 2
- [11] John H Flavell. Metacognition and cognitive monitoring: A new area of cognitive–developmental inquiry. *American psychologist*, 34(10):906, 1979. 1, 7, 9
- [12] Gregor Geigle, Radu Timofte, and Goran Glavaš. African or european swallow? benchmarking large vision-language models for fine-grained object classification. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2653–2669, 2024. 7
- [13] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. 2
- [14] Dedre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive science*, 7(2):155–170, 1983. 1, 2, 7, 9
- [15] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14375–14385, 2024. 7, 8
- [16] Stevan Harnad. The symbol grounding problem. *Physica D: Nonlinear Phenomena*, 42(1-3):335–346, 1990. 2
- [17] Ray S Jackendoff. *Languages of the mind: Essays on mental representation*. mit Press, 1995. 7
- [18] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Liuhui Wang, Jianhan Jin, Claire Guo, Shen Yan, et al. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning quality, robustness, and efficiency. *arXiv preprint arXiv:2502.09621*, 2025. 7, 8
- [19] Siddharth Joshi, Haoli Yin, Rishabh Adiga, Ricardo Monti, Aldo Carranza, Alex Fang, Alvin Deng, Amro Abbas, Brett Larsen, Cody Blakeney, et al. Datbench: Discriminative, faithful, and efficient vlm evaluations. *arXiv preprint arXiv:2601.02316*, 2026. 2
- [20] Jihyung Kil, Zheda Mai, Justin Lee, Arpita Chowdhury, Zihe Wang, Kerrie Cheng, Lemeng Wang, Ye Liu, and Wei-Lun Harry Chao. Mllm-compbench: A comparative reasoning benchmark for multimodal llms. *Advances in Neural Information Processing Systems*, 37:28798–28827, 2024. 7
- [21] Jeonghwan Kim and Heng Ji. Finer: Investigating and enhancing fine-grained visual concept recognition in large vision language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6187–6207, 2024. 7
- [22] Michael J Kolen and Robert L Brennan. *Linking. In Test equating, scaling, and Linking: methods and Practices*, pages 487–536. Springer, 2014. 14
- [23] Asher Koriat. How do we know that we know? the accessibility model of the feeling of knowing. *Psychological review*, 100(4):609, 1993. 7
- [24] Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. Naturalbench: Evaluating vision-language models on natural adversarial samples. *Advances in Neural Information Processing Systems*, 37:17044–17068, 2024. 7
- [25] Yijiang Li, Qingying Gao, Tianwei Zhao, Bingyang Wang, Haoran Sun, Haiyun Lyu, Robert D Hawkins, Nuno Vasconcelos, Tal Golan, Dezhi Luo, et al. Core knowledge deficits in multi-modal language models. *arXiv preprint arXiv:2410.10855*, 2024. 7, 8
- [26] Yian Li, Wentao Tian, Yang Jiao, Tianwen Qian, Na Zhao, Bin Zhu, Jingjing Chen, and Yu-Gang Jiang. Look before you decide: Prompting active deduction of mllms for assumptive reasoning. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 2713–2722, 2025. 8
- [27] Zhiyuan Li, Heng Wang, Dongnan Liu, Chaoyi Zhang, Ao Ma, Jieting Long, and Weidong Cai. Multimodal causal reasoning benchmark: Challenging vision large language models to discern causal links across modalities. *arXiv preprint arXiv:2408.08105*, 2024. 8

- 446 [28] John M Linacre et al. What do infit and outfit, mean-square
447 and standardized mean. *Rasch measurement transactions*, 16
448 (2):878, 2002. 14
- 449 [29] Fanqing Meng, Jin Wang, Chuanhao Li, Quanfeng Lu, Hao
450 Tian, Jiaqi Liao, Xizhou Zhu, Jifeng Dai, Yu Qiao, Ping
451 Luo, et al. Mmiu: Multimodal multi-image understanding
452 for evaluating large vision-language models. *arXiv preprint*
453 *arXiv:2408.02718*, 2024. 7
- 454 [30] Thomas Mensink, Jasper Uijlings, Lluís Castrejon, Arushi
455 Goel, Felipe Cadar, Howard Zhou, Fei Sha, André Araujo,
456 and Vittorio Ferrari. Encyclopedic vqa: Visual questions
457 about detailed properties of fine-grained categories. In *Pro-*
458 *ceedings of the IEEE/CVF International Conference on Com-*
459 *puter Vision*, pages 3113–3124, 2023. 7, 8
- 460 [31] Gregory Murphy. *The big book of concepts*. MIT press, 2004.
461 1, 7
- 462 [32] Gregory L Murphy and Douglas L Medin. The role of theories
463 in conceptual coherence. *Psychological review*, 92(3):289,
464 1985. 7
- 465 [33] Judea Pearl. *Causality*. Cambridge university press, 2009. 2,
466 7, 8, 9
- 467 [34] Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai
468 Sun, Gongjun Xu, and Mikhail Yurochkin. tinybench-
469 marks: evaluating llms with fewer examples. *arXiv preprint*
470 *arXiv:2402.14992*, 2024. 14
- 471 [35] Chenhui Qiang, Zhaoyang Wei, Xumeng Han, Zipeng Wang,
472 Siyao Li, Xiangyuan Lan, Jianbin Jiao, and Zhenjun Han.
473 Ver-bench: Evaluating mllms on reasoning with fine-grained
474 visual evidence. In *Proceedings of the 33rd ACM Interna-*
475 *tional Conference on Multimedia*, pages 12698–12705, 2025.
476 7, 8
- 477 [36] Georg Rasch. *Probabilistic models for some intelligence and*
478 *attainment tests*. ERIC, 1993. 12
- 479 [37] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-
480 Manor. Imagenet-21k pretraining for the masses. *arXiv*
481 *preprint arXiv:2104.10972*, 2021. 10, 19
- 482 [38] Eleanor Rosch. Principles of categorization. In *Cognition*
483 *and categorization*, pages 27–48. Routledge, 2024. 7
- 484 [39] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan
485 Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot:
486 Advancing multi-modal language models with a comprehen-
487 sive dataset and benchmark for chain-of-thought reasoning.
488 *Advances in Neural Information Processing Systems*, 37:8612–
489 8642, 2024. 7, 8
- 490 [40] Roger N Shepard. Toward a universal law of generalization
491 for psychological science. *Science*, 237(4820):1317–1323,
492 1987. 2
- 493 [41] Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen,
494 Kurt Mehlhorn, and Karsten M Borgwardt. Weisfeiler-lehman
495 graph kernels. *Journal of Machine Learning Research*, 12(9),
496 2011. 15
- 497 [42] Zechen Sun, Yisheng Xiao, Juntao Li, Yixin Ji, Wenliang
498 Chen, and Min Zhang. Exploring and mitigating shortcut
499 learning for generative large language models. In *Proceed-*
500 *ings of the 2024 joint international conference on computa-*
501 *tional linguistics, language resources and evaluation (LREC-*
502 *COLING 2024)*, pages 6883–6893, 2024. 2
- [43] Anne M Treisman and Garry Gelade. A feature-integration
theory of attention. *Cognitive psychology*, 12(1):97–136,
1980. 2, 7
- [44] Amos Tversky. Features of similarity. *Psychological review*,
84(4):327, 1977. 2, 7, 9
- [45] Lorraine K Tyler and Helen E Moss. Towards a distributed ac-
count of conceptual knowledge. *Trends in cognitive sciences*,
5(6):244–252, 2001. 1, 7
- [46] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui,
Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and
Serge Belongie. The inaturalist species classification and
detection dataset-supplementary material. *Reptilia*, 32(400):
1–3, 2017. 10
- [47] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun,
Tong Xu, and Enhong Chen. A survey on multimodal large
language models. *National Science Review*, 11(12):nwae403,
2024. 7
- [48] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai
Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang,
Huan Sun, et al. Mmmu-pro: A more robust multi-discipline
multimodal understanding benchmark. In *Proceedings of the*
63rd Annual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers), pages 15134–15186,
2025. 7
- [49] Zhiwei Zha, Xiangru Zhu, Yuanyi Xu, Chenghua Huang,
Jingping Liu, Zhixu Li, Xuwu Wang, Yanghua Xiao, Bei
Yang, and Xiaoxiao Xu. Constructure: Benchmarking concept
structure reasoning for multimodal large language models. In
Findings of the Association for Computational Linguistics:
EMNLP 2024, pages 4954–4968, 2024. 7
- [50] Jusheng Zhang, Kaitong Cai, Xiaoyang Guo, Sidi Liu, Qinhan
Lv, Ruiqi Chen, Jing Yang, Yijia Fan, Xiaofei Sun, Jian
Wang, et al. Mm-cot: a benchmark for probing visual chain-
of-thought reasoning in multimodal models. *arXiv preprint*
arXiv:2512.08228, 2025. 8
- [51] Letian Zhang, Xiaotong Zhai, Zhongkai Zhao, Yongshuo
Zong, Xin Wen, and Bingchen Zhao. What if the tv was
off? examining counterfactual reasoning abilities of multi-
modal language models. In *Proceedings of the IEEE/CVF*
conference on computer vision and pattern recognition, pages
21853–21862, 2024. 8

544 A. Towards A Unified Relational Framework

545 The pursuit of increasingly general artificial intelligence
546 foregrounds a foundational question: how should founda-
547 tional models or agents abstract and organize experience
548 into reusable conceptual knowledge [9]? Cognitive science
549 grounds this intuition: experience-distilled conceptual com-
550 ponents underpin basic cognitive operations and, in turn,
551 scaffold the higher-level capabilities required for real-world
552 tasks. Across decades of debate, cognitive science has con-
553 verged on one foundational principle: **conceptual knowl-**
554 **edge is internally organized and relational** [14, 32, 38].
555 In this paper, we formalize *multimodal conceptual structure*
556 (*MCS*) in multimodal LLMs, motivated by this convergent
557 view, as the grounded relational organization that determines
558 how a concept is internally constituted, how it is situated
559 among interconnected concepts, and how it can be causally
560 explained [17, 32].

561 **From Flat Features to Structured Concepts.** Classical
562 definitional theories treated concepts as sets of necessary-
563 and-sufficient features: an object either satisfies the checklist
564 or it does not. Rosch’s *prototype theory* overturned this view
565 by demonstrating that category membership is graded, *i.e.*,
566 some members are more typical than others, establishing
567 that concepts possess internal organization beyond binary
568 inclusion [38]. *Structure-mapping theory* shifted the focus
569 from the features themselves to the relations among them:
570 conceptual competence requires preserving structured corre-
571 spondences that support comparison, contrast, and analogi-
572 cal transfer [14]. *Theory-theory* pushed further still, arguing
573 that even relational similarity is insufficient without explana-
574 tory constraints: what lends a concept its coherence is an
575 underlying explanatory theory that specifies why features
576 co-occur and how they are causally linked [32]. On this
577 view, knowing that birds have wings and fly is not merely a
578 statistical regularity but reflects an implicit causal model in
579 which wings enable flight.

580 B. Extended Related Work

581 As summarized in Table 2, recent MLLM benchmarks have
582 become increasingly comprehensive, spanning vision-centric
583 evaluation [5, 24, 48], fine-grained recognition [12, 21],
584 web/knowledge base (KB) grounding [6, 30], comparative
585 and multi-image reasoning [20, 29, 49], and reasoning-trace
586 diagnostics [15, 18, 25, 35, 39]. However, these efforts are
587 rarely organized around a cognitively motivated ability lad-
588 der that explicitly distinguishes and diagnoses (i) perceptual
589 grounding and attribute binding [1, 43], (ii) relational infer-
590 ence over conceptual structure [14, 44], (iii) counterfactual
591 and interventional predictions [3, 33], and (iv) metacognitive
592 verification of reasoning integrity [11, 23]. This diagnos-
593 tic need is practically consequential to achieve reliable and

robust multimodal intelligent systems [47]. MCSBench is
designed to fill this gap by proposing a unified relational
benchmark to diagnose concept–attribute, concept–concept,
multi-image, and hypothetical reasoning in cognitive concep-
tual space, paired with integrity-aware evaluation. Complete
related works are in Appendix.

We expand the main-paper positioning (Table 2) along
three axes that jointly define the gap MCSBench is designed
to fill.

Recent MLLM benchmarks have made substantial
progress along complementary directions: vision-centric
evaluation, fine-grained recognition, web/knowledge base
(KB) grounding, comparative and multi-image reasoning,
and reasoning-trace diagnostics [5, 6, 12, 15, 18, 20, 21,
24, 25, 29, 30, 35, 39, 48, 49]. Two structural gaps, how-
ever, persist across this landscape. First, existing bench-
marks are rarely organized around a **cognitively motivated**
framework that explicitly distinguishes and diagnoses per-
ceptual grounding, relational inference, hypothetical reason-
ing, and metacognitive verification as progressively deeper
human cognitive operations [1, 3, 11, 44]. Second, much
concept/category-centric evaluation remains largely *entity-*
or label-centric, *i.e.*, testing whether a model can align a
concept name to visual evidence, retrieve associated facts,
or classify among semantic labels, rather than adopting a
relational view in which a concept is characterized by its
discriminative attributes, its conceptually confusing neigh-
bors, and the hypothetical transformations that changes its
identity/population [14, 31, 45]. Below, we discuss how
each gap emerges across three evaluation families.

B.1. From Entity-Centric to Relational Evaluation

Fine-grained benchmarks such as FOCI [12] and Finer [21]
aim to evaluate frontier MLLMs on discriminative seman-
tic cues among confusing categories, but their evaluation
unit remains the individual concept-label pair success means
aligning the correct label to visual evidence, without sys-
tematically probing the relational structure (attribute bind-
ings, neighborhood contrasts, hypothetical transformations)
that characterizes human conceptual organization [1, 31].
Comparative benchmarks extend the scope toward relational
reasoning, but along narrower axes: Constructure [49] evalu-
ates concept abstraction and specialization along taxonomic
hierarchies (*i.e.*, superordinate and subordinate), yet does not
probe lateral relations among confusers at arbitrary semantic
levels; MLLM-CompBench [20] and MMIU [29] benchmark
pairwise image relativity and multi-image understanding at
scale, yet neither controls for concept-level confusability
nor grounds its comparisons in the discriminative attribute
structure of a concept family. Noticeably, our primary focus
is probing internal conceptual alignment; we additionally
include an open-book retrieval setting as a diagnostic lens
on evidence utilization.

Benchmark	Task Design						Cognitive Abilities				
	Fine-gr.	Vision-centric	Ext-KB/Web	Comp./Relat.	Ctf./Interv.	Integrity eval.	Percep.	Concept	Relat.	Ctf.	Meta-verify
<i>General vision benchmarks</i>											
MMStar	✗	✓	✗	✗	✗	~	✓	✗	✗	✗	✗
MMMU-Pro	✗	✓	✗	~	✗	✗	✓	✓	~	✗	✗
NaturalBench	~	✓	✗	✓	✗	✗	✓	~	✓	✗	✗
<i>Fine-grained concept benchmarks</i>											
FOCI	✓	~	✗	✗	✗	✗	✓	✓	✗	✗	✗
Finer	✓	~	✗	✗	✗	✗	✓	✓	✗	✗	✗
<i>Knowledge-intensive / web-grounded VQA</i>											
Encyclopedic-VQA	✓	~	✓	✗	✗	~	✓	✓	~	✗	✗
InfoSeek	✗	~	✓	✗	✗	✗	✓	✓	~	✗	✗
<i>Concept-structure / relational reasoning</i>											
CONSTRUCTURE	~	~	✗	✓	✗	✗	✓	✓	✓	✗	✗
MLLM-CompBench	✗	~	✗	✓	✗	✗	✓	~	✓	✗	✗
MMIU	✗	✓	✗	✓	✗	~	✓	~	✓	✗	✗
<i>Evidence / reasoning-trace / integrity diagnostics</i>											
Visual CoT	~	✓	✗	✗	✗	~	✓	~	~	✗	~
VER-Bench	✓	✓	✗	✗	✗	~	✓	~	~	✗	~
MME-CoT	✗	~	✗	~	✗	✓	~	~	~	✗	✓
HallusionBench	✗	✓	✗	~	✗	✓	~	~	~	✗	✓
CoreCognition	✗	~	✗	~	~	✓	~	✓	~	~	✓
MCSBench (Ours)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 2. Comparison of MCSBench with representative benchmarks across task-design axes (fine-grained, vision-centric, external knowledge, comparative/relational, counterfactual, integrity evaluation) and cognitive ability dimensions (perception, concept knowledge, relational and counterfactual reasoning, metacognitive verification). ✓: explicitly present as evaluation objective; ~: implicitly present; ✗: not targeted.

646 MCSBench adopts the relational cognitive view explicitly:
 647 every task is constructed around a local conceptual-space
 648 topology, *i.e.*, a target concept situated among semantically
 649 close confusers sharing controlled attribute overlap, and re-
 650 quires specific relational operations within that topology,
 651 progressing from attribute–concept binding (L1) through
 652 inter-concept contrast and cross-image matching (L2) to
 653 hypothetical transformation under attribute edits (L3) and
 654 metacognitive verification of reasoning integrity (L4). Suc-
 655 cess therefore requires reasoning over the structure that sys-
 656 tematically defines a concept in the complex network, not
 657 merely recognizing or labeling it in isolation.

658 B.2. Hypothetical Reasoning in MLLMs

659 Hypothetical and counterfactual reasoning has received
 660 growing attention as a diagnostic dimension beyond recogni-
 661 tion. C-VQA [51] augments VQA questions with counter-
 662 factual presuppositions, revealing sharp performance drops,
 663 yet its counterfactuals operate at the scene level and are dis-
 664 tinct from conceptual structure; CFMM [26] and MuCR [27]
 665 probe broad causal or counterfactual competence through
 666 binary classification and siamese image pairs, yet neither
 667 tests how attribute-level interventions shift concept identity
 668 within a controlled confuser set; MM-CoT [50] evaluates
 669 chain-of-thought under adversarial distractors but remains
 670 verification-centric rather than structure-centric. MCSBench
 671 introduces hypothetical reasoning over the conceptual topol-
 672 ogy, an advanced prediction capability of human [3, 33].

673 B.3. Evidence and Reasoning-Trace Diagnostics

674 Prior work evaluates reasoning traces and evidence uti-
 675 lization along complementary but structurally unanchored

676 axes: Visual CoT [39], VER-Bench [35], and MME-
 677 CoT [18] assess chain-of-thought quality across percep-
 678 tion and cognition subtasks; HallusionBench [15] and
 679 CoreCognition [25] diagnose hallucination-style failures;
 680 Encyclopedic-VQA [30] and InfoSeek [6] measure evi-
 681 dence utilization yet equate success with answer correctness
 682 alone. Since none is organized around a controlled concept-
 683 centric topology, attributing failures to specific conceptual
 684 structural deficits remains difficult.

685 C. Additional Details on Dataset Profile

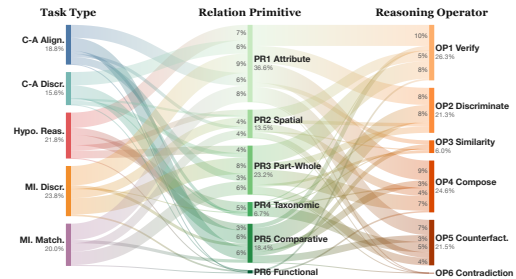


Figure 1. Alluvial diagram of MCSBench composition across three axes: task type (left), relation primitive (center), and reasoning operator (right). Edge width encodes the proportion of items linking each pair; percentages denote the marginal share of each node.

686 **MCSBench Task Taxonomy** We define seven base tasks
 687 spanning four conceptual reasoning categories and four cog-
 688 nitive levels. Let q denote the query task, c a concept, I_c the
 689 object-centric image of c , R the relation set, A the attribute
 690 set, and $T(\cdot)$ a templated prompt.

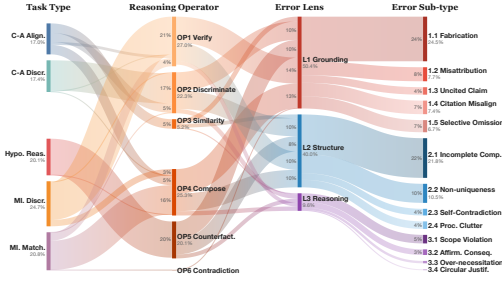


Figure 2. Error distribution of MCSBench overlay across task type, reasoning operator, error lens, and error sub-type.

691 *Concept-Attribute Alignment* (C-A Align.) aims to correctly
692 align the attribute relations with the query concept encoded
693 in the image. Given an image and task query, the model
694 predicts the correct concept relation $\Pr(A | T(q), I_c)$.

695 *Concept-Attribute Discrimination* (CA-Discr.) asks the
696 model to pick the correct attribution a^+ over a distractor
697 a^- for the same question, i.e. $\Pr(A^+ \succ A^- | T(q), I_c)$,
698 distinguishing the concept from its nearest confusable neigh-
699 bor [44].

700 *Concept-Concept Alignment* (CC-Align.) requires following
701 a specified relation R to identify the correct compositional
702 answer. Given $T(q)$, R , and I_c , the model prefers the cor-
703 rect candidate: $\Pr(A^+ \succ A^- | T(q), R, I_c)$. Because
704 candidates come from the same taxonomic neighborhood,
705 success depends on tracking the *relational role* of R , not
706 category-level similarity [14].

707 *Concept-Concept Discrimination* (CC-Discr.) checks
708 whether the model can jointly recover the structural order-
709 ing σ , relation R , and concept c' from the query and image:
710 $\Pr(R, c' | T(q), I_c)$. Unlike CC-Align, no relation label is
711 given, so the model must infer both the relation type and its
712 filler while rejecting spurious same-neighborhood pairings.

713 *Hypo. Reas. (Interventional)* evaluates how the model
714 responds to an explicit intervention on an attribute. Fol-
715 lowing Pearl’s do-calculus [33], given $T(q, \text{do}(A \leftarrow$
716 $A'))$ and I_c , the model predicts the concept $\Pr(c' |$
717 $T(q, \text{do}(A \leftarrow A')), I_c)$. Failure suggests the model relies
718 on co-occurrence rather than relational structure.

719 *Hypo. Reas. (Counterfactual)* focuses on counterfactual rea-
720 soning for a fully observed factual instance. Given evidence
721 $e := (T(q), \tau, R, I_c)$ that fixes the realized instance, the
722 model predicts the counterfactual concept under the alter-
723 native attribute value A' : $\Pr(c_{A \leftarrow A'} = c^{\text{cf}} | T(q), R, I_c)$,
724 where $c_{A \leftarrow A'}$ is the potential-outcome variable defined by
725 Pearl’s abduction–action–prediction.

726 *Multi-Image Discrimination* (MI-Discr.) asks the model
727 to identify which images in a candidate set satisfy a textual
728 description of attributes or relations. The output is a selection
729 (or ranking) of images, evaluated by whether true positives

are preferred over visually similar distractors, so the task
 requires grounding relational specs across multiple images
 at once.

Multi-Image Matching (MI-Match.) tests whether a model
 can produce a globally consistent correspondence between
 multiple image instances and multiple textual descriptions.
 Given a set of images and a set of descriptions, the model
 outputs a multi-to-multi (or constrained) assignment between
 them.

Reasoning-Chain Evaluation assesses whether a model’s
 proposed explanation is structurally valid w.r.t. the question
 and the provided evidence. We construct candidates where
 the final answer can be correct even when intermediate steps
 are unsupported or logically inconsistent, thereby measuring
 the SAI overlay [11] gap between answer correctness and
 justified correctness.

Fine-grained Perspective on Structure. Fig. 3(c) decom-
 poses each item along two axes: the *relation primitive* it
 engages—spanning attribute/morphological (PR1), spatial
 (PR2), part-whole (PR3), taxonomic (PR4), and compar-
 ative (PR5) relations—and the *reasoning operator* it de-
 mands, ranging from verification (OP1) and discrimination
 (OP2) through compositional constraint satisfaction (OP4)
 to counterfactual intervention (OP5). The distribution is
 intentionally weighted toward PR1 (34.3%) and PR3 (26.4%),
 emphasizing perceptually grounded relational distinctions
 over coarse taxonomic knowledge. Crucially, the many-to-
 many flow confirms that no task type reduces to a single
 relation–operator pair: Concept–Attribute Alignment routes
 primarily through OP1 and OP4, while Hypothetical Reason-
 ing concentrates on OP5 despite drawing on equally diverse
 primitives. This ensures that cross-task performance differ-
 ences reflect the demanded cognitive operation rather than
 surface content shifts.

Relation-Operator Decomposition. Fig. 1: We themati-
 cally annotate each reasoning-chain question along two axes:
 relation primitives (PR1–PR6, multi-label), i.e., the seman-
 tic relation engaged by the question body and options, and
 a primary reasoning operator (OP1–OP6), i.e., the cogni-
 tive move required to answer correctly. The distribution of
 MCSBench is intentionally weighted toward perceptually
 grounded relations with PR1 Attribute (36.6%), PR3 Part-
 Whole (23.2%), and PR5 Comparative (18.4%) over coarse
 taxonomic knowledge (PR4, 6.7%). Operators spread more
 evenly across verification (OP1, 26.3%), composition (OP4,
 24.6%), counterfactual (OP5, 21.5%), and discrimination
 (OP2, 21.3%). Crucially, routing is many-to-many: each
 task type draws on distinct primitive–operator combinations,
 confirming that cross-task performance gaps reflect the de-
 manded cognitive operation rather than surface content over-
 lap.

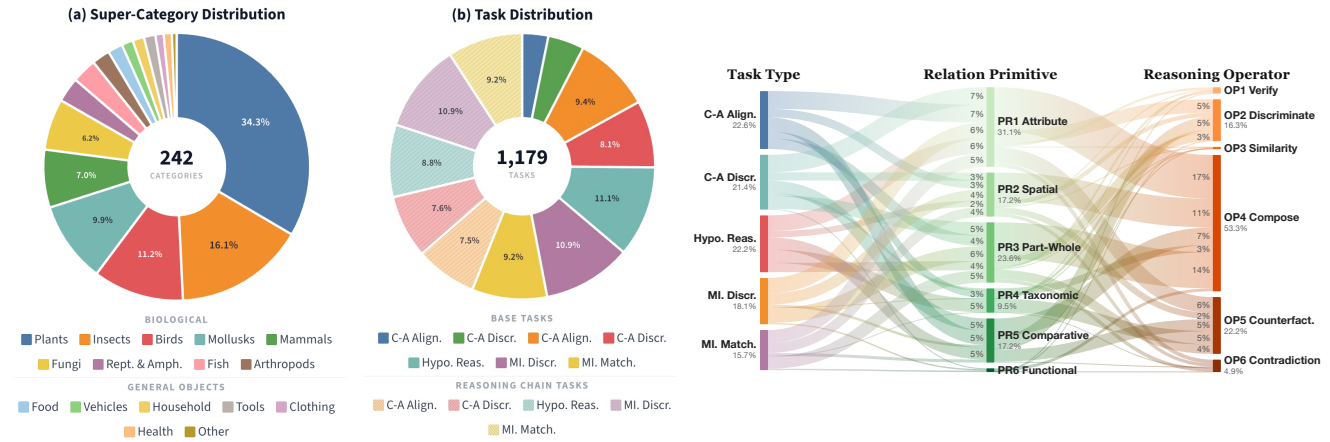


Figure 3. Our MCSBench is profiled with (a) Supercategory distribution of 242 categories. (b) 661 questions in 7 base tasks, where 518 are paired with reasoning tasks. (c) Fine-grained structural primitives.

781 **Structural Error Taxonomy.** We annotate each distractor in the reasoning-chain overlay with a primary error drawn from a three-lens hierarchy applied bottom-up (Fig. 2).
 782 L1 Grounding errors (50.4%) concern misalignment between claims and the provided evidence or observations:
 783 fabrication of unsupported attributes (24.5%), misattribution of a correct trait to the wrong entity (7.7%), citation misalignment where a cited source does not support the claim (7.4%), and selective omission of required constraints (6.7%).
 784 L2 Structure errors (40.0%) concern argument completeness: incomplete comparison, *i.e.*, confirming one option without eliminating plausible confusers (21.8%), and non-uniqueness, where the stated cues remain compatible with multiple candidates (10.5%).
 785 L3 Reasoning errors (9.6%) involve formal logical faults such as scope violations (*e.g.*, using out-of-scope criteria (*e.g.*, habitat) as decisive evidence (5.0%)) and affirming the consequent (3.0%).
 786 The steep concentration in L1/L2 reveals that the primary bottleneck is evidentiary grounding and contrastive elimination which are precisely the structural competencies our task targets.

801 D. Details on Dataset Construction & Quality Control

802 This section specifies our integral implementations on constructing MCSBench with quality control.

805 D.1. Dataset Construction

806 **Building Concept Bank.** MCSBench evaluates conceptual structure centered on multimodal concepts sourced from a general-purpose concept bank construction pipeline. Each concept is documented with its linguistic name, image set, definition, meronymic attributes, and a set of visually confusing distractors. We source concept names from three visually rich taxonomies: iNat2021 [46] (10,000 fine-grained species across 11 supercategories), ImageNet [8] (1,000 ob-

814 ject categories spanning everyday and abstract concepts), and
 815 ImageNet-21K [37] (the *small* subset of $\sim 10,000$ synsets sampled from the full 21,841-class hierarchy). Noticeably,
 816 we only adopt ImageNet images in the benchmark: our manual inspection reveals that a substantial fraction of iNat2021
 817 images fail to unambiguously depict the target species, often proving indistinguishable among visually confusing distractors or even absent. We therefore retain iNat2021 solely as a concept-name source. Concept names are canonicalized and deduplicated. Intangible concepts (15 rejection super-categories: abstract concepts, emotions, processes, events, relationships, etc.) are also filtered. The curated list contains 15,343 unique and concrete visual concepts. *Next*, we construct a meronymic hierarchy (part-whole attributes) for each concept by prompting a web-augmented agent (o4-mini) based on authoritative sources with Prompt (details see Section I). The obtained hierarchical attributes act as contextual references for the task generation prompts. *Lastly*, we derive visually confusing concepts from a foundational vision encoder (Open-CLIP [7]) and linguistically confusing concepts from web-augmented agents. An optional expansion stage further enriches diversity: additional confusing concepts are iteratively generated and added to the bank, with corresponding images crawled from public sources (*e.g.*, WikiCommons). Noisy images are automatically filtered by subsequent pipeline stages.

840 **Generating Base Tasks.** Base items are formatted as 5-way multiple-choice VQA (one correct option). The implementation uses task-specific prompt conditioned generators that collectively instantiate the seven base tasks in Table 2.

841 **Prompting and distractor control.** For each sampled (concept, image(s)) pair, the generator is conditioned on concept name/definition, meronymic attribute hierarchy, and (when applicable) the confusable concept set. Prompts of

848 each generation is additionally contextualized with the task-
849 specific distractor strategy (e.g., visual similarity, linguistic
850 proximity), enabling both controlled hardness and avoiding
851 shortcuts. In our runs, a web-augmented LLM produces 5
852 questions per image, each with 5 options; outputs are format-
853 validated (key presence, option format correctness, answer
854 index validity), with automatic retries on malformed outputs;
855 unanswerable images are discarded.

856 **Sampling scale (generation pool).** To balance diversity and
857 cost, we sample 100 concepts per concept source partition
858 for confusing-category tasks and 50 concepts per partition
859 for single-image and multi-image tasks, with up to 3 images
860 per concept. This yields a raw pool of thousands of candidate
861 questions before filtering, which is progressively curated by
862 the quality-control pipeline into the final benchmark.

863 **Committee-Based Refinement (Auditor–Challenger–
864 Testee).** We implement the committee in Fig. 2 as an
865 iterative refinement loop that targets three failure modes:

- 866 • **Challenger (hardening):** refines or regenerates items to
867 prevent trivial elimination strategies (e.g., distractors that
868 are syntactically inconsistent, or apparently solvable with-
869 out visual evidence), and to increase distractor plausibility
870 while preserving the target test relation.
- 871 • **Auditor (faithfulness):** checks overall quality of the ques-
872 tion in four aspects: Factual Correctness (answer verifi-
873 able), Distractor Correctness (distractor demonstrably in-
874 correct), Image Dependency (visual information involved),
875 and Self-Containedness (no missing context or ambiguity);
876 it then filters likely hallucinated questions based on the
877 quality score, minimal reference number (≥ 3), and target
878 attribute presence.
- 879 • **Testee (difficulty sanity):** probes whether items are de-
880 generate (too easy via priors or has modality shortcut);
881 trivial items are discarded.

882 In our implementation, this committee behavior is instan-
883 tiated through (i) optional challenge refinement iterations
884 (notably for single-image tasks), and (ii) staged automated
885 auditing used as gating signals: quality scoring, hallucina-
886 tion/grounding checks, answer/distractor verification, and
887 justification sufficiency (detailed in D.2).

888 **Constructing Reasoning Chains (Integrity Overlay).** For
889 base tasks with complex relations ($\geq L2$), we construct
890 reasoning-chain evaluation items that probe whether correct-
891 ness is supported by structurally valid reasoning (measured
892 by SAI). For each retained base item, an Explainer produces
893 a gold reasoning chain (with controlled quality) grounded in
894 (i) explicit visual observations and (ii) atomic textual claims
895 referencing authoritative sources. A red-team attacker then
896 generates structurally perturbed distractors, plausible but in-
897 correct chains with diverse error types (see Fig. 2). One base
898 task is paired with three RC questions. Each RC question
899 is assembled as a multiple-choice selection over candidate
900 chains by sampling 4 combinations from 5 generated dis-

tractors and inserting the gold chain at a random position. 901
In total, the benchmark contains 661 base items spanning 7 902
tasks. Of these, 518 contain the integrity overlay ($\geq L2$), pro- 903
ducing 1,554 reasoning-chain evaluation items (3 per base 904
item). 905

906 D.2. Quality Control

907 **Overview.** Our quality control is a four-stage pipeline:
908 (1) shortcut filtering removes items solvable from language
909 alone; (2) LLM judgment applies automated multi-criteria
910 adjudication for correctness/grounding/plausibility; (3) hu-
911 man curation removes residual ambiguity and validates opti-
912 mality; (4) optional psychometric filtering calibrates item
913 difficulty and screens potentially biased items via measure-
914 ment modeling.

915 **Shortcut Filtering.** We detect language-only shortcuts by
916 evaluating each candidate item under text-only settings us-
917 ing a model ensemble (Gemini-2.5-Pro, Gemini-2.5-Flash,
918 o4mini as Testee) and discarding items that remain solvable
919 without images. Specifically, we run a shortcut detector that
920 flags an item if text-only accuracy exceeds a fixed thresh-
921 old, indicating likely reliance on linguistic priors rather than
922 image grounding.

923 **LLM Judgment** for multi-criteria check. Surviving items
924 are filtered by an automated LLM adjudication stack that
925 evaluates complementary failure modes:

- 926 1. **Quality check (Auditor):** assigns a 1–5 quality score and
927 sub-scores over factual correctness, distractor validity,
928 image dependency, and self-containedness.
- 929 2. **Hallucination check:** removes items whose claimed evi-
930 dence is ungrounded, inconsistent, or not aligned with the
931 image and references, using stricter criteria for settings
932 prone to hidden priors.
- 933 3. **Option verification:** verifies that only the answer option
934 is correct and optimal, and that each distractor is plausible
935 yet incorrect; items failing answerability or with high
936 leakage are discarded or refined.
- 937 4. **Justification (for RC overlay):** evaluates whether the
938 gold reasoning chain with referenced evidence is suffi-
939 cient, factually correct, and logically sound and coherent,
940 and filters out items where the reasoning cannot be reli-
941 ably justified given the evidence.

942 **Human Curation** After automated filtering, we conduct two
943 rounds of human curation, including annotation and consis-
944 tency check, on Prolific platform with increased stringency
945 to confirm that (i) the gold option is correct and optimal,
946 (ii) distractors are plausible and incorrect, and (iii)
947 each question is answerable and clear. Among all partic-
948 ipants, all native English speakers residing in the United
949 States (60.8%) or the United Kingdom (39.2%) and holding
950 postgraduate degrees (81.0% Master’s, 19.0% Doctoral), en-
951 suring the reading comprehension and domain knowledge
952 required for nuanced visual-reasoning judgments. The pool

953 was demographically balanced (59.5% female; mean age
954 37.2, $\sigma=8.1$, range 23–50; five self-reported ethnic groups)
955 and experienced, with a median of 1,457 previously ap-
956 proved Prolific tasks, reducing annotator bias and strength-
957 ening ground-truth reliability.

958 **Psychometric Filtering.** As the final quality-control stage,
959 we optionally calibrate items with a Rasch IPL model so
960 that the released benchmark functions as a stable measure-
961 ment instrument. We retain items that fall within a target
962 difficulty range ($[-1.5, 1.5]$ logits) and satisfy item-fit and
963 discrimination screens, and then stratify the retained set into
964 easy, medium, and hard tiers for controlled evaluation. Full
965 model specification, filtering criteria, sampling procedure,
966 and bootstrap estimation details are provided in App. F.4.

967 E. Details on Evaluation Protocol

968 **Evaluation protocol.** All models are evaluated through
969 a unified inference pipeline with greedy decoding
970 (`do_sample=False`, $t=0.0$, `max_new_tokens=4096`)
971 with fixed random seed to ensure deterministic generation.
972 The default instruction asks the model to select from options
973 A–E (or F–J for reasoning-chain evaluation) and end with the
974 phrase without enforcing or suppressing chain-of-thought;
975 this lets each model respond in its natural style while en-
976 suring fairness. Images are inserted before the text prompt
977 following each model’s native multi-image template, resized
978 to a maximum dimension of 448 px with the same aspect-
979 ratio. Answers are parsed by a multi-stage rule-based parser
980 along with GPT LLM with the highest scores retained. For
981 API-served models, we use $t=0.0$ (whenever feasible) with
982 asynchronous inference (`concurrency=30`) and exponential-
983 backoff retries; all other settings are identical. Non-API
984 models are inferred locally with $8 \times$ A6000 NVIDIA GPUs.
985 The evaluation scripts and prompts will be released with the
986 benchmark.

987 F. Details on Experiments & Analysis

988 F.1. Main Results

989 F.2. Psychometrics Benchmarking via IRT

990 A diagnostic benchmark should function as a reliable mea-
991 surement instrument. We apply Rasch IPL to a task-
992 balanced variance-sampled short form (149 items) to val-
993 idate MCSBench under both ACC and SAI scoring, and to
994 quantify how SAI reshapes model ordering. Model abilities
995 are estimated with 95% bootstrap CIs ($B=200$); SAI abili-
996 ties are linked to the ACC scale via mean–sigma linking for
997 comparison. Full details are in the Appendix.

998 **Reliable measurement.** Both scoring schemes achieve high
999 person-separation reliability (ACC: 0.973; SAI: 0.980), with
1000 all items passing infit acceptance (Table 8), confirming MCS-
1001 Bench stably separates model abilities.

Table 5. **Per-task sensitivity to semantic intervention.** For each task, we report the mean accuracy change Δ (pp) across $n=12$ models, paired t -test statistic t , and effect size d (Cohen’s d). * $p < 0.05$.

Task	Δ (pp)	t	p	d
C-A Align.	-1.16	-0.77	0.459	-0.22
C-A Discr.	+5.88	2.74	0.019*	0.79
C-C Align.	+4.28	2.53	0.028*	0.73
C-C Discr.	+1.22	1.43	0.180	0.41
Hypo. Reas.	+0.89	1.26	0.235	0.36
MI-Discr.	-2.99	-2.01	0.070	-0.58
MI-Match.	-0.08	-0.08	0.934	-0.02
Overall	+0.87	2.88	0.015*	0.83

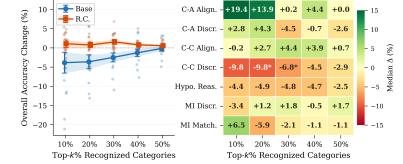


Figure 4. **Recognition-conditioned evaluation.** MCSBench is re-evaluated on the subset of the top- k % recognized categories ($k \in \{10, 20, 30, 40, 50\}$). **Left:** Overall accuracy trend for Base and reasoning-chain (R.C.) tasks across models. **Right:** Per-task median accuracy change Δ (pp) under the same filtering; *significant under Wilcoxon signed-rank test.

1002 **SAI imposes stricter difficulty.** Mean item difficulty shifts
1003 from $b=0.038$ (ACC) to $b=0.477$ (SAI), as SAI penalizes
1004 “lucky-correct” responses lacking structural justification.

1005 **Preserved global order, non-trivial local re-ranking.**
1006 Model ordering is largely maintained (Kendall $\tau=0.870$)
1007 with no systematic directional bias ($p=0.525$, n.s.), yet re-
1008 rankings are substantial (mean $|\Delta\text{Rank}|=2.80$, max = 10).
1009 Table 9 shows that models such as Kimi-VL-A3B drop 10
1010 positions under SAI while Gemma-3-12B and MiniCPM-V-
1011 4 each rise by 7, exposing structural gaps that ACC obscures.
1012 SAI thus serves as a *conservative structural refinement*: pre-
1013 serving broad rank structure while revealing where model
1014 correctness is structurally unjustified.

1015 F.3. Recognition-Conditioned Evaluation

1016 **Setup.** Since open-ended generation is difficult to reliably
1017 evaluate across MLLMs, we operationalize recognition as
1018 a multiple-choice classification task. For each category, we
1019 randomly sample 10 images and present the model with
1020 a forced-choice question identifying the depicted concept.
1021 Distractors are constructed by retrieving the nearest neigh-
1022 bors in OpenCLIP embedding space, ensuring semantically
1023 proximate options that prevent trivial elimination.

1024 **Additional Results.** We present additional performance sen-
1025 sitivity for base and reasoning tasks.

1026 F.4. Item-Response Modeling

1027 **Implementation details.** To characterize model behavioral
1028 patterns and calibrate measurement, we leverage a Rasch
1029 IPL model [36], which places model ability θ_i and item dif-
1030 ficulty b_j on a common logit scale: $P(X_{ij}=1 | \theta_i, b_j) =$
1031 $\sigma(\theta_i - b_j)$. Difficulties are estimated by marginal maximum
1032 likelihood [2] (girth, 61-point Gauss–Legendre quadra-
1033 ture, $\mathcal{N}(0, 1)$ prior) and abilities by expected *a posteriori*
1034 inference. We calibrate under two views: binary accuracy
1035 and SAI for both correctness and integrity. Starting from 518

Table 3. Performance evaluation on MCSBench base tasks under the closed-book setup. *Align.=Alignment; Discr.=Discrimination; Match.=Matching; Cog.=Cognitive; Con.=Concept; Attr.=Attribute; Hypo.=Hypothetical.*

Model	Con.-Attr.		Con.-Con.		Hypo.	Multi-Image		Cog. Ladder			Avg.
	Align.	Discr.	Align.	Discr.	Reason.	Discr.	Match.	L0	L1	L2	
<i>Proprietary Models</i>											
GPT-5	91.67	84.31	89.58	61.07	64.86	85.16	87.96	87.36	81.72	61.07	78.37
Claude-Opus-4.6	91.67	82.35	83.33	53.44	64.86	84.38	85.19	86.21	79.46	53.44	75.19
Gemini-2.5-Pro	88.89	82.35	81.25	50.38	61.26	87.50	85.19	85.06	79.01	50.38	74.13
Gemini-2.5-Flash	77.78	76.47	73.96	37.40	55.86	78.12	84.26	77.01	73.14	37.40	66.57
O4-Mini	77.78	70.59	81.25	35.11	43.24	85.94	85.19	73.56	74.04	35.11	66.26
GPT-5-Mini	69.44	68.63	76.04	34.35	47.75	81.25	79.63	68.97	71.33	34.35	63.69
Claude-4.5-Sonnet	91.67	78.43	75.00	32.06	37.84	72.66	62.96	83.91	62.08	32.06	59.00
<i>Open-Sourced Models</i>											
InternVL3-78B-Instruct	83.33	78.43	77.08	25.19	38.74	72.66	77.78	80.46	66.37	25.19	60.06
Qwen2.5VL-72B-Instruct	80.56	74.51	72.92	28.24	39.64	77.34	74.07	77.01	66.14	28.24	60.06
Qwen3VL-30B-A3B-Instruct	77.78	58.82	69.79	27.48	43.24	71.88	81.48	66.67	66.59	27.48	58.85
GLM-4.5V	80.56	70.59	72.92	22.90	36.94	74.22	76.85	55.17	60.95	20.61	58.09
InternVL3.5-38B-Instruct	80.56	64.71	69.79	24.43	31.53	70.31	68.52	71.26	60.05	24.43	54.46
MiniCPM-V-4.5	63.89	37.25	57.29	21.37	31.53	75.78	64.81	48.28	58.01	21.37	49.47
Ovis2.5-9B	77.78	54.90	61.46	21.37	31.53	57.81	63.89	64.37	53.50	21.37	48.56
Gemma-3-27B-it	69.44	47.06	58.33	16.79	27.03	64.84	60.19	56.32	52.82	16.79	46.14
Kimi-VL-A3B-Instruct	75.00	47.06	67.71	16.79	26.13	57.81	57.41	58.62	51.92	16.79	45.84
NVILA-15B	77.78	60.78	69.79	24.43	33.33	43.75	35.19	67.82	44.70	24.43	43.72
DeepSeek-VL2	44.44	37.25	60.42	26.72	24.32	59.38	49.07	40.23	48.31	26.72	42.97
Phi-4-Multimodal-Instruct	63.89	37.25	59.38	19.08	14.41	50.78	41.67	48.28	41.31	19.08	37.82
LLaVA-OV-1.5-8B-Instruct	88.89	64.71	64.58	18.32	25.23	27.34	27.78	74.71	34.99	18.32	36.91
R1-Onevision-7B	52.78	33.33	48.96	20.61	23.42	46.88	35.19	41.38	38.60	20.61	35.40
Llama3.2-11B-Vision-Instruct	61.11	41.18	64.58	17.56	23.42	14.84	15.74	49.43	27.99	17.56	28.74
Idefics3-8B-LLama3	52.78	43.14	55.21	17.56	18.02	19.53	21.30	47.13	27.31	17.56	27.99
LongVILA-R1-7B	55.56	47.06	51.04	17.56	25.23	8.59	25.93	50.57	26.19	17.56	27.69
<i>Thinking Models</i>											
Qwen3VL-30B-A3B-Thinking	83.33	66.67	64.58	32.06	32.43	81.25	76.85	73.56	64.33	32.06	59.15
GLM-4.1V-9B-Thinking	75.00	58.82	67.71	25.95	23.42	68.75	68.52	65.52	57.11	25.95	52.04
Kimi-VL-A3B-Thinking-2506	69.44	50.98	53.12	18.32	26.13	53.91	59.26	47.13	41.53	17.56	43.57

Table 4. Performance on MCSBENCH with Reasoning-Chain evaluation overlaid across five base tasks. We report base accuracy and SAI per task, averaged across all questions. Open-book settings provide models with web search (*web*) or ground-truth evidence (*evd*) as context. Models are ranked by descending average SAI within each group.

Model	C-C Discr.		C-C Align.		Hyp. Reas.		MI-Discr.		MI-Match.		Avg.	
	Base	SAI	Base	SAI	Base	SAI	Base	SAI	Base	SAI	Base	SAI
<i>Open-Book Eval</i>												
Gemini-3-Flash-Preview (web)	84.44	82.59	76.14	68.94	74.04	70.83	80.47	78.91	84.26	80.25	79.92	76.51
Gemini-3-Pro-Preview (web)	85.56	79.26	79.55	71.21	79.81	72.44	86.72	82.55	87.04	75.00	83.98	76.45
GPT-5.2-2025-12-11 (web)	80.00	74.81	61.36	50.38	39.42	35.26	86.72	78.12	83.33	68.83	71.04	62.29
GPT-5-2025-08-07 (web)	87.78	70.00	68.18	53.41	58.65	41.99	89.06	75.52	84.26	64.20	78.19	61.71
InternVL3-78B-Instruct (evd)	77.78	77.04	40.91	38.26	24.04	22.44	72.66	71.35	77.78	74.69	59.46	57.59
Qwen2.5-VL-72B-Instruct (evd)	74.44	68.52	42.05	35.23	27.88	24.04	77.34	76.56	74.07	70.99	60.23	56.44
Qwen3-VL-30B-A3B-Instruct (evd)	67.78	60.37	44.32	39.39	27.88	23.72	71.88	67.19	81.48	71.60	59.65	53.47
MiniCPM-V-4.5 (evd)	56.67	52.96	30.68	23.86	21.15	16.67	75.78	71.88	64.81	52.47	51.54	45.30
<i>Proprietary Models—Closed-Book Eval</i>												
Gemini-3-Pro-Preview	84.44	77.41	73.86	65.53	72.12	66.35	90.62	85.42	89.81	75.31	82.82	74.71
Claude-Opus-4.6	82.22	75.93	68.18	56.44	51.92	49.04	84.38	79.69	85.19	88.58	75.93	68.15
O4-Mini	80.00	68.15	45.45	37.12	36.54	30.45	85.94	75.00	85.19	70.68	67.95	57.53
GPT-5-2025-08-07	88.89	62.22	64.77	48.11	58.65	39.10	85.16	66.41	87.96	61.11	77.61	55.98
GPT-5.2-2025-12-11	70.00	63.70	42.05	34.47	39.42	31.41	68.75	59.90	53.70	45.68	55.41	47.55
<i>Open-Sourced Models—Closed-Book Eval</i>												
InternVL3-38B-Instruct	78.12	75.00	33.33	26.13	25.19	19.85	65.62	59.38	79.63	62.96	58.55	50.23
Qwen2.5-VL-72B-Instruct	74.44	61.85	42.05	35.98	27.88	24.36	77.34	64.84	74.07	55.56	60.23	49.36
Qwen3-VL-30B-A3B-Instruct	67.78	47.04	44.32	31.06	27.88	22.44	71.88	52.60	81.48	59.26	59.65	43.31
Gemma-3-27B-it	55.56	33.70	25.00	18.56	19.23	14.42	64.84	49.48	60.19	37.96	46.33	32.05
Kimi-VL-A3B-Thinking-2506	45.56	18.52	30.68	14.02	21.15	11.54	38.28	16.67	55.56	10.80	38.42	14.29
LLaVA-OneVision-1.5-8B-Instruct	62.22	23.33	26.14	11.74	21.15	5.13	27.34	10.94	27.78	10.49	32.05	11.97

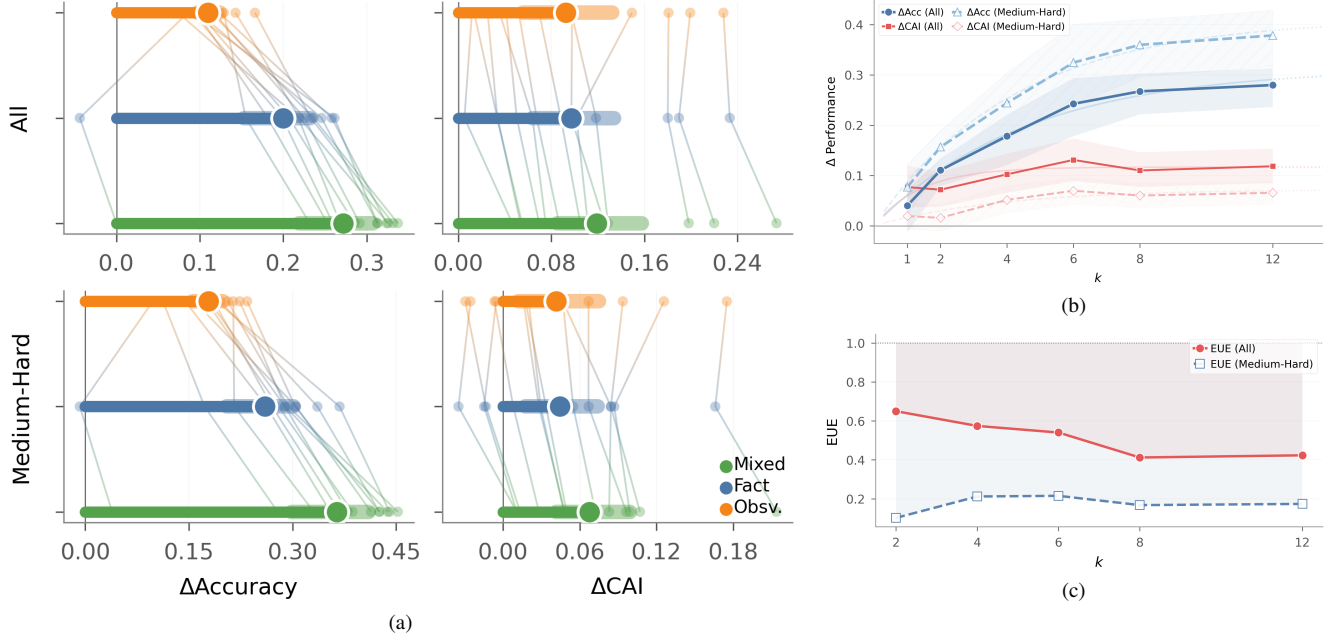


Figure 5. **In-context diagnostics: evidence type, scaling, and utility.** (a) Model-wise Δ Accuracy and Δ SAI gains over the no-evidence baseline across three evidence types, for ALL and MEDIUM-HARD subsets. (b) Δ Accuracy and Δ SAI trends w.r.t. random evidence count k ; Acc gains saturate; SAI gains minorly. (c) Evidence-Utilization Efficiency (EUE) across k ; $EUE < 1$ indicates low evidence gains on reasoning integrity.

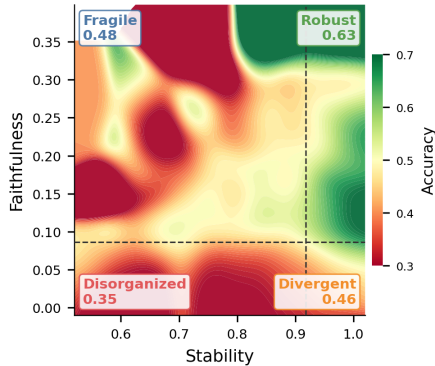


Figure 6. Reasoning Stability–Faithfulness landscape. Contours show binned mean accuracy (w/ thinking rollout); dashed lines mark median thresholds partitioning items into four regimes (Robust, Fragile, Divergent, Disorganized) with inset mean accuracies. $N = 2310$ model \times question pairs.

1036 items across 56 models, to satisfy Rasch model’s assumption, a two-stage screen [28] removes 180 misfitting items
 1037 (infit/outfit $MNSQ \notin [0.5, 1.5]$) and 21 low-discriminating
 1038 items (point-biserial $r < 0.15$); per fit iteration, variance-
 1039 based proportional sampling [34] then yields a 149-item
 1040 short form. SAI abilities are linked to the ACC scale via
 1041 mean–sigma equating [22] ($A=0.7497$, $B=0.0974$), with
 1042 uncertainty from bootstrap resampling ($N=200$, percentile
 1043 95% CIs). Both views achieve person-separation reliabil-
 1044 ity ≥ 0.97 , confirming stable measurement while exposing
 1045

Table 6. Graph-structure diagnostics for MCSBench. (a) Spearman ρ between Stability and accuracy across Easy/Hard split (Fisher’s z : $p < .05$). (b) Top OLS predictors of accuracy on reasoning graph features.

(a) Difficulty interaction (Spearman ρ)			
Metric	ρ_{Easy}	ρ_{Hard}	Fisher z
Stability	+0.29	−0.14	8.59*
(b) Top structural predictors (OLS, std. β)			
Feature	β	p	Rank
Edge Type	+0.099	<.001	1
Stability	+0.049	<.001	2
Faithfulness	+0.042	<.001	3

ranking shifts that accuracy alone conceals. 1046

Additional Results. We present entire model ability ranking 1047
 and delta rank in Fig. 12 9. The distribution of capability 1048
 and local rank adjustment validate our claims in the main 1049
 paper. 1050

F.5. Reasoning Graph Analysis 1051

For each model \times question pair we sample $k=5$ stochastic 1052
 chain-of-thought rollouts (temperature $t \in \{0.7, 0.8\}$) and 1053
 extract typed reasoning graphs from the resulting rationales 1054
 using Gemini-3-Pro-Preview with a fixed structured-JSON 1055

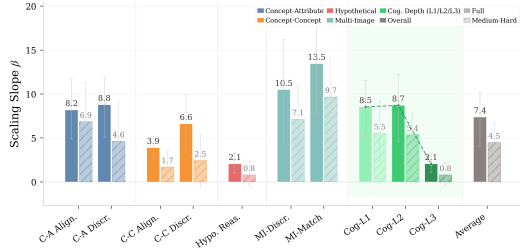


Figure 7. **Scaling slopes β per task** (OLS on $\ln P$) with 95% bootstrap CIs, coloured by task group. Scale alone accounts for moderate variance ($R_{sc}^2 \approx .32$); family intercepts contribute the largest incremental gain ($\Delta R_{fam}^2 \approx .50$), indicating that recipe-level shifts dominate over scaling rates.

Table 7. Variance decomposition of task-level accuracy via nested OLS. R_{scale}^2 : $\ln P$ alone; ΔR_{fam}^2 : family intercepts; R_{full}^2 : full model. * $p < .05$; $\dagger p < .10$.

Task	R_{scale}^2	ΔR_{fam}^2	R_{full}^2
C-A Align.	.209	.418*	.691
C-A Discr.	.320	.287*	.747
C-C Align.	.393	.320*	.820
C-C Discr.	.252	.288*	.661
Hypo-Reas.	.179	.285 \dagger	.642
MI-Discr.	.200	.618*	.866
MI-Match.	.300	.513*	.867
Avg.	.321	.500*	.870

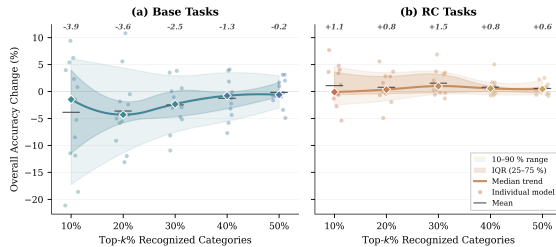


Figure 8. Performance sensitivity to removing top- k % most recognized categories. Base accuracy even minorly degrades, while RC accuracy remains flat, indicating that recognition does not aid correctness and reasoning integrity. Bands: IQR and 10–90% range across models.

1056 schema. Each graph comprises six node types (concept, 1057
img, obs, fact, claim, answer) and five edge 1058
types (grounds, supports, refutes, compares, 1059
concludes), assigned deterministically via keyword- 1060
based rules. An edge is *evidence-anchored* if it is a 1061
grounds edge originating from a visual-evidence node 1062
(img/obs). The full extraction pipeline is released with the 1063
benchmark.

Graph similarity kernel. We quantify structural similarity with a Weisfeiler–Leman (WL) subtree kernel [41] using $h=3$ iterations, node-type label initialization, directed neighborhoods encoded with OUT/IN markers and edge-type signatures. Two summary scores are derived from this kernel:

- *Stability*: mean pairwise WL similarity across all $\binom{5}{2}=10$ rollout pairs per question, capturing whether a model converges on a single reasoning strategy or vacillates across stochastic samples.
- *Faithfulness*: mean WL similarity between each rollout graph and the ground-truth reference graph extracted from the released golden reasoning chains (`rc_breakdown`, available for 518/661 items), measuring how closely the model’s spontaneous reasoning structure mirrors the intended solution path, independent of answer correctness.

F.6. Scaling Trend

Implementation details. We evaluate 63 proprietary models or open-source VLMs spanning 1.98B–107.71B parameters, representing scale as $x = \log_{10}(\text{params}_B)$. Each model is scored on MCSBENCH under two splits: FULL (661 items) and MEDIUM+HARD (319 items), reporting both ACC and SAI at the task and cognitive-depth level.

Scaling estimation. Scaling slopes are fitted via ordinary least squares (OLS) on the log-linear form $y = A + \beta x + \epsilon$ for each task and cognitive depth group, with 95% confidence intervals obtained by bootstrap resampling over models ($N_{boot}=2000$, $\alpha=0.05$).

Scaling trend visualization. We visualize the per-task scaling trend in Fig. 10. We can observe that the variance is larger among the medium-sized models, which further suggests that other confounders, *e.g.*, training recipe, have larger effects, supporting our main claims.

G. Additional Results

H. External Consistency

To establish external consistency, we correlate MCSBench overall accuracy with three established benchmarks, *i.e.*, MME, HallusionBench, and SEEDBench, on overlapping model sets (setup in Appendix). Table ?? reports statistically significant rank correlations across all three suites, confirming that MCSBench tracks the core visual–language competence axis in the community. Fig. ?? further reveals that MCSBench tasks provide a diagnostically richer decomposition. Most notably, SAI uniquely correlates with Hallucination Robustness, confirming that structural-awareness recovers robustness invisible to raw accuracy. We also notice that *Hypo. Reas.* has a lower correlation, suggesting that it may capture missing capabilities.

Table 8. IRT calibration (Rasch IPL) comparing ACC and SAI scoring on 149 variance-sampled items (95% bootstrap CIs, $B = 200$). Reliability: Rasch person-separation reliability; mean b : item difficulty (logits). Rank agreement between scoring schemes is evaluated after linking SAI to ACC scale.

Key Stats	ACC	SAI
Reliability	0.973 $^{+.002}_{-.001}$	0.980 $^{+.001}_{-.001}$
Infit acceptable	100%	100%
Mean b (logits)	0.038	0.477
Kendall τ	0.870 $^{+.020}_{-.033}$	
Mean $ \Delta Rank $	2.80 $^{+.66}_{-.37}$ (Max = 10)	
Wilcoxon (Rank)	$p = 0.525$ (<i>n.s.</i>)	

Table 9. Largest rank movers when switching from ACC to structure-aware SAI, sorted by $|\Delta R|$. $\Delta\theta^* = \theta_{SAI}^* - \theta_{ACC}$ is the linked IRT ability shift. \blacktriangle rank improves; \blacktriangledown rank worsens. Full ranking in Appendix.

Model	ΔR	$\Delta\theta^*$
Kimi-VL-A3B	\blacktriangledown 10	-0.481
Gemma-3-12B	\blacktriangle 7	+0.388
MiniCPM-V-4	\blacktriangle 7	+0.396
InternVL3.5-20B-A4B	\blacktriangle 6	+0.384
DeepSeek-VL2-Small	\blacktriangledown 5	-0.298
LLaVA-OV-Qwen2-7B	\blacktriangledown 5	-0.641
GPT-5	\blacktriangledown 3	-0.698
Qwen2.5-VL-72B	\blacktriangle 3	+0.158
Gemini-3-Flash	\blacktriangle 2	+0.014

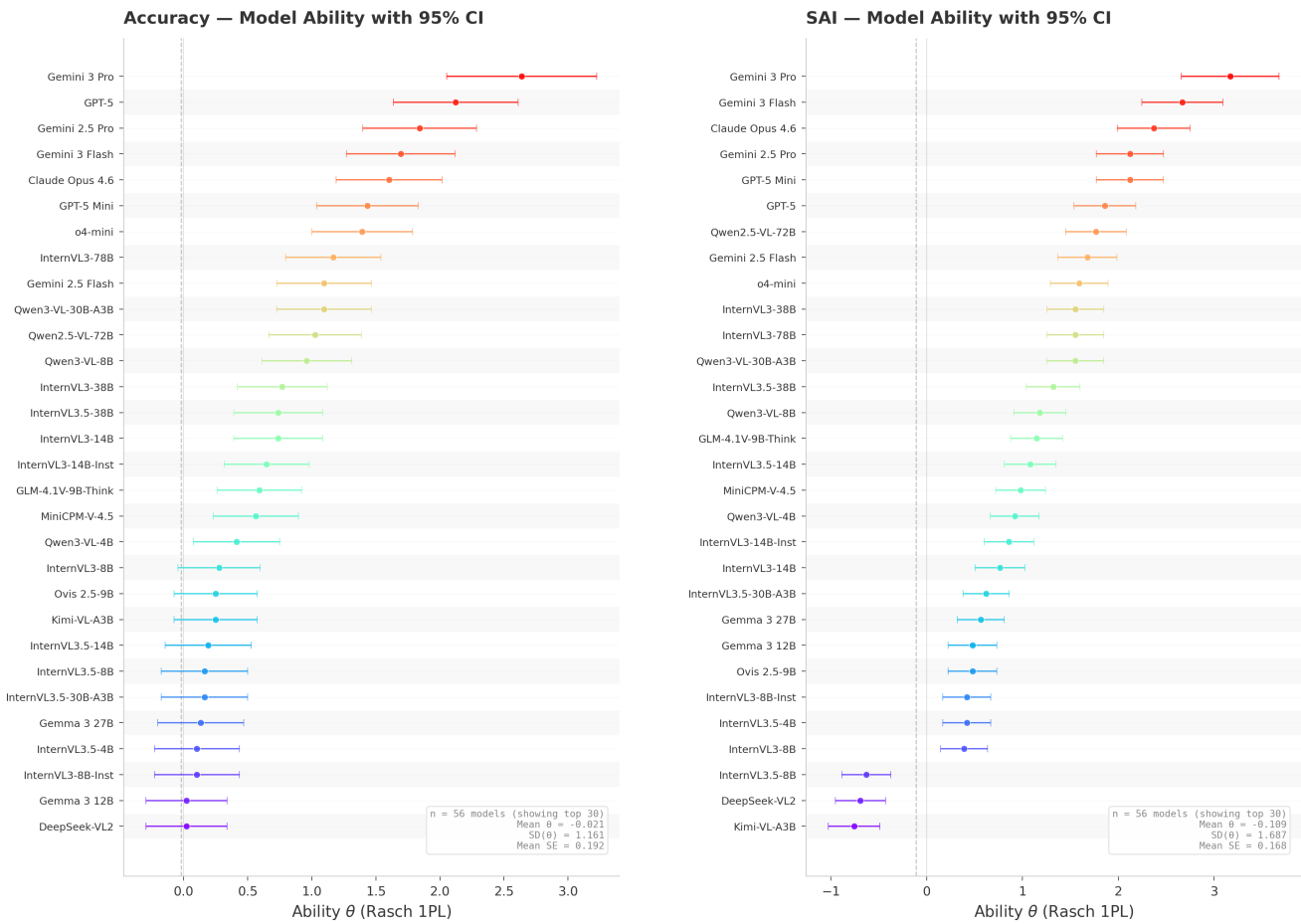


Figure 9. Rasch model ability estimates for top-30 models ranked by accuracy (left) and SAI (right).

1111 **H.1. Prompt Sensitivity Analysis**

1112 Five prompting variants, Direct, Chain-of-Thought (CoT),
 1113 Justification, Expert Role, and Contrast, are compared across
 1114 9 models with complete prompt grids (Fig. 11). Four conclu-

sions emerge: (1) overall accuracy is only mildly sensitive to
 prompt wording; (2) CoT prompting does not improve perform-
 ance over simple direct instruction; (3) Contrast prompting
 is consistently harmful, suggesting that current models lack
 robust discriminative comparison ability; and (4) within each

1115
 1116
 1117
 1118
 1119

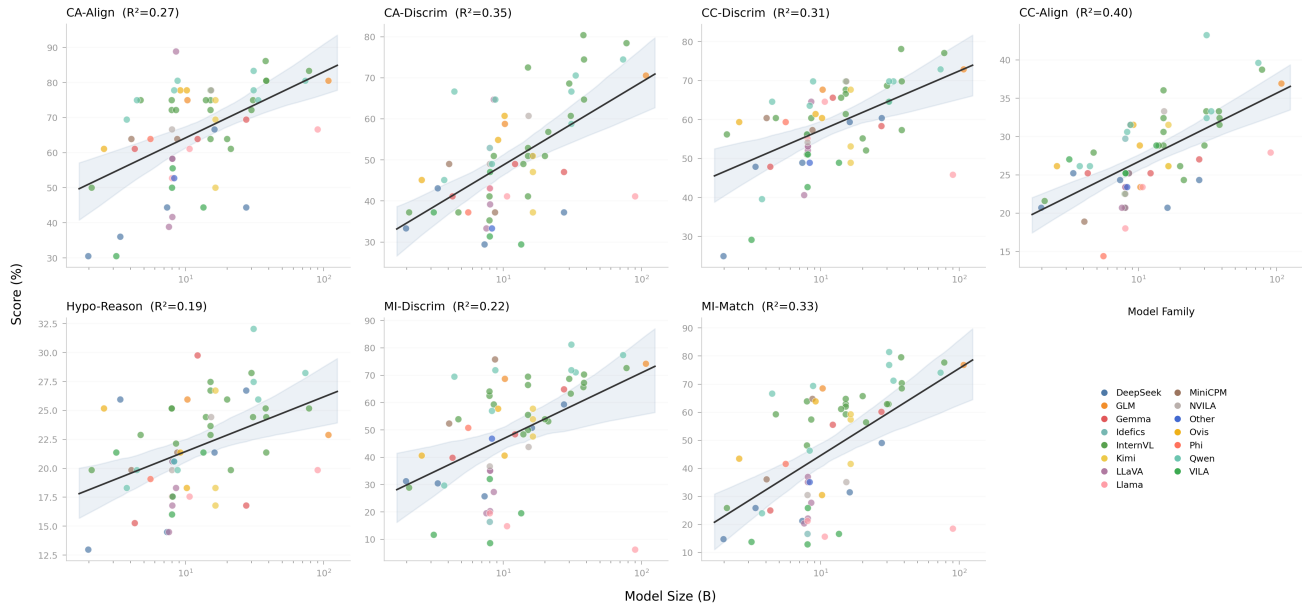


Figure 10. Per-task scaling trend visualization.

1120 cognitive level, the difficulty ordering is preserved across
1121 all prompts—prompt engineering modulates absolute scores
1122 but does not reshape the benchmark’s difficulty profile.

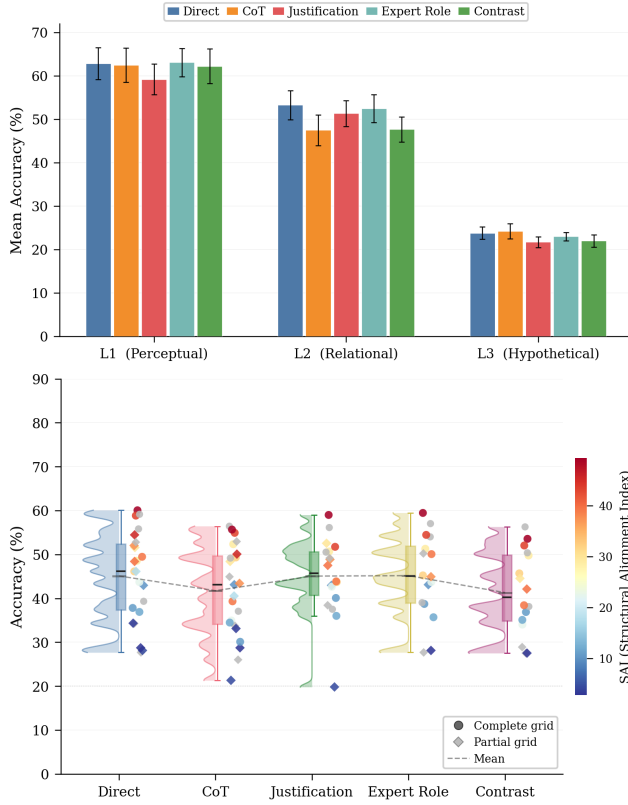


Figure 11. **Prompt sensitivity across cognitive levels and prompt type.** *Left:* Mean accuracy under five prompt variants aggregated over 9 models, with error bars showing between-model variability. Prompt effects are largest at L2 and negligible at L3. *Right:* Per-prompt accuracy distributions; marker color encodes SAI and shape denotes grid completeness. No prompt yields a uniform gain.

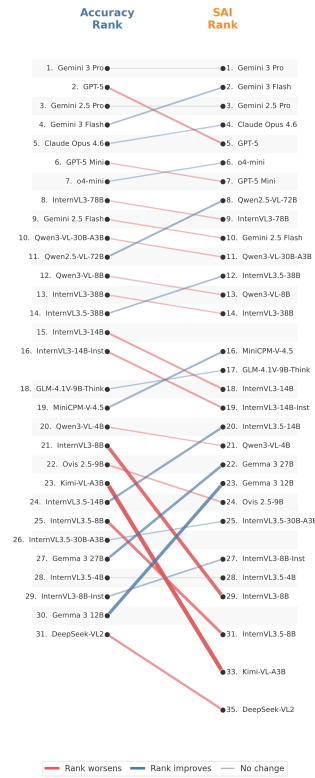


Figure 12. IRT-calibrated rank shift from accuracy to SAI. Rasch ability estimates reorder the leaderboard substantially, confirming that correctness alone does not reflect reasoning integrity.

1123 I. Prompts

1124 I.1. Definition: Relation Primitives

1125 We annotate each question with one or more *relation primitives* that characterize the semantic relation(s) explicitly
1126 invoked by the question stem and answer options. Tags are
1127 **multi-label**: a single item may involve multiple primitives
1128 (e.g., an attribute comparison over a part of an object).
1129

- 1130 • **PR1 Attribute / Morphological.** Visual attributes and
1131 intrinsic properties such as color, shape, texture, size, ma-
1132 terial, or state (*has-attribute*, *has-value*).
- 1133 • **PR2 Spatial Configuration.** Relative position and layout
1134 relations such as above/below, inside/outside, next-to, ori-
1135 entation, or arrangement (*spatially-related-to*,
1136 *relative-position*).
- 1137 • **PR3 Mereological / Part-Whole.** Part structure and
1138 compositional relations including subcomponents, con-
1139 tainment, and component membership (*has-part*,
1140 *part-of*, *component-of*).
- 1141 • **PR4 Taxonomic / Definitional.** Category membership
1142 and definitional relations such as class/subclass, instance-
1143 of, or definitional inclusion (*is-a*, *defined-as*,
1144 *instance-of*).
- 1145 • **PR5 Comparative / Contrastive.** Explicit compar-
1146 isons or contrastive judgments (e.g., *more/less*,
1147 *same/different*, *similarity/contrast*) over concepts or at-
1148 tributes (*differs-from*, *more/less-than*).
- 1149 • **PR6 Functional / Contextual.** Purpose, causal, or contex-
1150 tual/ecological relations (e.g., *used for*, *caused by*, *occurs*
1151 *in a situation or environment*) (*used-for*, *caused-by*,
1152 *occurs-in*).
- 1153 • **PR0 Miscellaneous.** The relation is ambiguous, under-
1154 specified, or does not fit PR1–PR6.

1155 **Annotation rule.** Tagging is determined solely by the re-
1156 lations *explicitly* engaged by the question and options; an-
1157 notators do *not* infer additional relations from background
1158 knowledge.

1159 I.2. Definition: Reasoning Operators

1160 Each question is assigned exactly one primary operator and
1161 up to one secondary operator. Selection follows a strict
1162 priority ordering: OP6 > OP5 > OP4 > OP2 > OP3 >
1163 OP1.

- 1164 • **OP1 Verify / Ground:** verify that a concept or image
1165 satisfies a described attribute profile.
- 1166 • **OP2 Discriminate:** identify the discriminative at-
1167 tribute(s) separating a target from confusers.
- 1168 • **OP3 Similarity / Align:** identify the closest concept
1169 under shared-attribute constraints.
- 1170 • **OP4 Compose / Bind:** satisfy multiple constraints si-
1171 multaneously; conjunction or multi-image matching.

- **OP5 Counterfactual / Intervene:** hypothetical at- 1172
tribute edit mapped to the nearest concept. 1173
- **OP6 Contradiction / Consistency:** locate which 1174
clause is contradicted by evidence or image. 1175

I.3. Relation Classification Prompt 1176

The following prompt is used to classify each question along 1177
both dimensions via an LLM-based annotation pipeline. 1178

```

ROLE: You are an expert conceptual-structure
annotator for a multimodal benchmark (MCSBench).
Given a question stem and its answer options,
classify the question along two dimensions: Relation
Primitives and Reasoning Operator.

INPUT:
Question: {ori.question}
Answer options: {ori.options}
Correct answer: {ori.answer}

OUTPUT FORMAT: Return only valid JSON:
{
  "rationale": "<1--2 sentence analysis>",
  "pr": ["PR#", ...],
  "op_primary": "OP#",
  "op_secondary": "OP#" | null
}

```

1179

I.4. Prompts for Dataset Generation 1180

J. Discussion and Limitations 1181

K. Data Licensing 1182

MCSBench contains source images from Wikimedia Com- 1183
mons (WikiCM)¹ and ImageNet-21K (IN21K) [8, 37]. All 1184
WikiCM images are dominantly verified to carry permissive 1185
Creative Commons licenses (CC0, CC BY, or CC BY-SA) at 1186
the time of collection. No images with *Non-Commercial* or 1187
NoDerivatives restrictions were included. Per-image license 1188
metadata and author attribution are provided in the released 1189
dataset files. IN21K images are governed by a research-only 1190
license that prohibits redistribution;². Accordingly, we re- 1191
lease only filenames and synset IDs so that authorised users 1192
may reconstruct the image set on their own. All MCSBench 1193
annotations will be released under CC BY-SA 4.0. which is 1194
intended for only non-commercial academic research. 1195

¹<https://commons.wikimedia.org>

²<https://image-net.org/accessagreement>



Figure 13. Multi-round quality control pipeline. Each part has conf./single-image/MI variants sharing the same structure.

Evidence-Based Explanation (conf. variant shown; single/MI analogous)**Role:** Domain expert in `{category_name}` and evidence-first adjudicator.**Goal:**

1. Extract every piece of nomenclature from the task body (question + options).
2. Explain each term in plain language; list common synonyms.
3. Provide authoritative web references for definitions and answer correctness.
4. Construct step-by-step justification grounded only in collected evidence.

Input: Image, confusing categories `{confusing_categories_list}`, VQA task (§4).**Output:** JSON with `attributes` (nomenclature terms), `references` (web sources), `evidence` (atomic facts from image/web), `reasoning` (step-by-step chain).**Justification Quality Filter** (conf. variant shown; MI analogous)**Role:** Domain expert in `{category_name}` and impartial dataset auditor.**Goal:**

1. Evaluate the quality of the reasoning chain produced by *Evidence-Based Explanation*.
2. If incorrect, *minimally tweak* the chain to make it correct.

Input: Image, confusing categories, original VQA task (§4), reasoning chain context (§5): Observations (O#), Evidence (E#), Reasoning steps.**Output:** Quality verdict + corrected chain if needed.**Reasoning-Chain (RC) Adversarial Attack** (conf. variant shown; single/MI analogous)**Role:** Reasoning-chain adjudicator and *adversarial task designer*. Has web search access.**Goal:** Given a VQA task + its reasoning chain for the correct answer, generate a *new* MC-VQA task where each option is an **incorrect reasoning chain for the same answer**.**Inputs:**

- Image of `{category_name}`; `{category_def}`
- Confusing categories: `{confusing_categories_list}`
- Original VQA task with answer (§7)
- REASON_CHAIN: JSON with `attributes`, `references`, `evidence`, `reasoning` (§8)

Procedure:

1. **Evaluate** the reasoning chain: factual accuracy, logical soundness, image grounding.
2. **Generate distractors:** 4 incorrect reasoning chains (options A–D) that are plausible but flawed.
3. Each distractor must contain a *subtle error* detectable only by an expert:
 - ▷ Fabricated observation, misattributed evidence, uncited claim
 - ▷ Incomplete comparison, self-contradiction, procedural clutter
 - ▷ Scope violation, affirming the consequent, over-necessitation, circular justification

Output: New MC task — question asks “Which reasoning chain correctly justifies the answer?”Options: 1 correct chain (original) + 4 adversarial chains. JSON with `error_type` per distractor.

Figure 14. Evidence-based explanation (5a) and justification quality filter (5b) produce and verify reasoning chains. RC adversarial attack generates MC tasks where distractors are *plausible but subtly flawed* reasoning chains.

Task Def – Identifying Discriminating Features (C-A Align., conf.)

Goal: Pinpoint the specific attribute(s) that best describes the *visual associations* (similarity, difference) between the image category T and confusing categories $[C]$.

Question Variations:

- ▷ “Which single attribute best distinguishes T from C [and $C_2, C_3 \dots$]?”
- ▷ “What is the *minimal set* of attributes needed to distinguish T from C_1 and C_2 ?”

Distractor Strategies:

1. Non-Minimal Set — includes redundant features
2. Insufficient Set — fails to distinguish from *all* specified C s
3. Non-Distinguishing (Shared) — T and all C s are visually similar on these attributes
4. Partially Discriminating — distinguishes from *some* but not all C s
5. Subtle/Quantitative — differs only in expert-level quantitative manner

Task Def – Mapping Attributes to Categories (C-A Discr., conf.)

Goal: Identify which confusing category C matches specified criteria for similarity/difference with target T on given attribute(s) $[A], [B] \dots$

Question Variations:

- ▷ “Which category is most visually similar to T regarding attribute A ?”
- ▷ “Which C is similar to T on A but different on B ?”

Distractor Strategies:

1. Opposite Relationship — opposite relation on the specified attribute
2. Mismatched Criteria — matches on different, unasked-for criteria
3. Partial Match — satisfies some but not all required criteria
4. Target Inclusion — T itself as an option (self-reference error)
5. Subtle Match — nuanced similarity that is hard to verify visually
6. Satisfies Canonically, Not Visually — correct for the species generally, wrong for *this image*

Task Def – Hypothetical Reasoning (conf.)

Goal: Predict which confusing category C a target T would most resemble if specific visual attributes were hypothetically modified (changed, added, removed).

Question Variations:

- ▷ “If T 's [Attr. A] were changed to [Value V], which C would the modified T most resemble?”
- ▷ “Which hypothetical modification would make T most similar to C ?”

Distractor Strategies:

1. Target as Distractor (T) — when the modification is significant
2. Similarity Decrease — C becomes *less* similar after modification
3. Persistent Difference — C remains clearly distinct even after modification
4. Second-Best Match — second most similar to the modified T
5. Close Relative of Correct — taxonomically close to correct C but subtly different

Figure 15. Task type definitions for **confusing-category** tasks (C-A Alignment, C-A Discrimination, Hypothetical Reasoning). Each definition specifies the goal, question variations, and task-specific distractor strategies injected into the generation prompt.

<p>Task Def – Image → Visual Description Matching (single-image)</p> <p>Goal: Present candidate attribute–description pairs; choose the one that most accurately describes a feature visible in the image.</p> <p>Distractor Strategies:</p> <ol style="list-style-type: none"> 1. Correct Attribute, Incorrect Description 2. Incorrect Attribute, Correct Description 3. Generalization Mismatch — plausible for the category but inaccurate for <i>this</i> image 4. Less Significant/Precise Match — less prominent feature than the correct answer 5. Conditional State Mismatch — correct if subject were in a different state/angle 6. Minor Feature Match — real but much less defining feature
<p>Task Def – Detecting Discrepancies (single-image)</p> <p>Goal: Present a plausible but incorrect statement describing the image. Identify which specific detail is contradicted by visible evidence.</p> <p>Distractor Strategies:</p> <ol style="list-style-type: none"> 1. Correct Detail Identification — detail that is actually <i>correct</i> 2. Non-Visual/Unverifiable Error — cannot be confirmed from the image alone 3. Minor Error Focus — trivial inaccuracy when a more significant error exists 4. Imprecise Error Description — incomplete description of the error <p>⋮ + 3 more strategies (<i>unrelated attribute, hypothetical context, less critical</i>)</p>
<p>Task Def – Image Selection by Attributes (multi-image)</p> <p>Goal: Identify a specific visual instance from a set of 4 closely related images, based on a textual description of 1–2 key diagnostic visual attributes. Category names are <i>not</i> revealed; options reference images directly (Image-1, ..., Image-4).</p> <p>Distractor Strategies:</p> <ol style="list-style-type: none"> 1. Confusing Attribute — shared attribute visually confusing among multiple images 2. Partial Match — matches one attribute but fails on another 3. Qualifier Distractor — shared attribute with different visual qualifiers/values 4. Indirect Reference — uses subtle visual descriptions instead of direct image references 5. Set Confusion — adds/removes one borderline image from the correct set

Figure 16. Task type definitions for **single-image** tasks (Image Matching, Discrepancy Detection) and **multi-image** tasks (Image Selection by Attributes). Together with Figure 15, these 7 task types cover the complete MCSBench taxonomy.

<p>Phase 0a: Supercategory Assignment</p> <p>Task: For each fine-grained category listed below, provide a list of relevant supercategory names.</p> <p>Requirements:</p> <ul style="list-style-type: none"> ▷ Supercategories should be common and not overly broad or vague. ▷ Provide at least 4 supercategory names; favor colloquial names. ▷ Rank from most specific (less coarse) to most general (coarser). <p>Input: <code>{fine-grained categories}</code></p> <p>Output: JSON — <code>{'concept': ['supcat_1', 'supcat_2', ...]}</code></p>
<p>Phase 0b: Concept Tangibility Filter</p> <p>Role: Filter WordNet synset names from ImageNet21K for tangible, visible objects.</p> <p>Accept if: concrete physical entity with defined boundaries, clearly photographable, consistent visual appearance, can be physically touched.</p> <p>Reject if falls into meta-categories: abstract concept, emotion, motion/action, scene without boundaries, temporal concept, process/phenomenon, event, relationship, property, social construct, theory, natural phenomenon, activity, ...</p> <p>Input: Batch of <code>{n_batch}</code> synset names with definitions.</p> <p>Output: JSON — <code>{'synset': 'ACCEPT/REJECT', 'meta-category': ..., 'explanation': ...}</code></p>

Figure 17. Phase 0: Concept Preparation. (a) Supercategory assignment builds the taxonomy hierarchy. (b) Tangibility filter removes abstract/non-visual concepts from the candidate pool.

Phase 1: MC-VQA Task Generation (confusing-category variant shown; single-image & multi-image variants analogous)

Role: Domain expert for `{category_name}`.
Goal: Design challenging, domain-specific MC questions testing fine-grained visual knowledge of `{category_name}`.

§1 Provided Information:

- Image of `{category_name}`; `{category_def}`
- Textual list of professional visual attributes (§9)
- List of visually confusing categories: `{confusing_categories_list}`

§2 Your Task:

1. Analyze image, attribute list, and confusing categories (internal only).
2. Generate `{num_questions}` MC tasks; craft an internal plan for difficulty.
3. Web search for additional information about confusing categories as needed.

§3 Question Content: Visual focus, evidence-based, domain attributes only. *Task Definition:* `{task}` (one of 7 task types from Figure 15)

§4 Choice Format: 5 choices (A–E); exactly 1 correct + 4 distractors.

§5 Strict Constraints:

- *Visual Verification* — all options must reference visible/verifiable features; image takes priority over text.
- *Professional Nomenclature* — formal domain-specific terms only; no colloquial hints.
- *Content Restrictions* — avoid revealing `{category_name}` or broader classification; no quantitative measurements.
- *Anti-Shortcut* — cannot be answered by elimination or text knowledge alone.

§6 Distractor Requirements: Visual similarity, terminology proximity, expert-level challenge, plausibility, category integration, task-specific challenges, knowledge–image conflict.

§7 Output: JSON list — `[{'question': ..., 'options': [A...E], 'answer': ..., 'strategy': ...}]`

§8 Exception Handling: Return error JSON if image blurred/irrelevant, attributes corrupted, etc.

§9 Attribute list for `{category_name}`: `{attribute_list}`

Figure 18. Phase 1: Main MC-VQA task generation prompt (confusing-category variant). The prompt injects a task-type definition (§3) from one of 7 task types. Single-image and multi-image variants share the same structure, differing only in §1 inputs and category references.

Phase 1 (Challenge): Iterative Difficulty Refinement (2 iterations; single/conf/MI variants)

Role: Expert biologist, evaluator and refiner of AI-generated visual assessments. *Has access to web search.*
Goal: Enhance difficulty of existing MC-VQA tasks so they are significantly more challenging for both expert humans and advanced MLLMs.

§1 Provided Information: Image of `{category_name}`, attribute list (§7), existing Q&A tasks (§8), web search, original task definition: `{task}`.

§2 Your Task:

1. Review image, attributes, and existing tasks; create internal mental description.
2. Plan to strategically increase difficulty following all requirements (§3–5).
3. Web search to connect with confusing categories, incorporate subtle visual info, ensure factual correctness.
4. Refine tasks: maximize difficulty by revising questions and options (especially distractors).

§3 Core Refinement Criteria:

- A. Maximize distractor plausibility & confusion.
- B. Apply task-specific strategies + broader misleading principles: `{challenge}`
- C. Prevent short-cuts: minimal redundancy, consistent format across options, require visual inspection.
- D. Maintain original task constraints and definition alignment.

§4 Constraints: (same as generation: visual verification, nomenclature, content restrictions)

Meta-Strategy Injected via `{challenge}`: *Subtle detail inaccuracies; attribute misattributions; plausible but absent features; incorrect spatial/relational descriptions; near-synonym substitutions; phylogenetic neighbor confusion; trait value shifts; contextual re-framing.*

Output: JSON with `''refinement_notes''` explaining changes (iteration 0: 4 choices; iteration 1: 5 choices).

Figure 19. Phase 1 (Challenge): Iterative difficulty refinement. Applied twice ($k=0, 1$) after initial generation. The meta-strategy list is injected via the `{challenge}` variable. Each iteration receives the output of the previous round as §8.

Phase 0a: Supercategory Assignment
<p>Task: For each fine-grained category listed below, provide a list of relevant supercategory names.</p> <p>Requirements:</p> <ul style="list-style-type: none"> ▷ Supercategories should be common and not overly broad or vague. ▷ Provide at least 4 supercategory names; favor colloquial names. ▷ Rank from most specific (less coarse) to most general (coarser). <p>Input: <code>{fine-grained categories}</code></p> <p>Output: JSON — <code>{`concept`: [`supcat_1`, `supcat_2`, ...]}</code></p>
Phase 0b: Concept Tangibility Filter
<p>Role: Filter WordNet synset names from ImageNet21K for tangible, visible objects.</p> <p>Accept if: concrete physical entity with defined boundaries, clearly photographable, consistent visual appearance, can be physically touched.</p> <p>Reject if falls into meta-categories: abstract concept, emotion, motion/action, scene without boundaries, temporal concept, process/phenomenon, event, relationship, property, social construct, theory, natural phenomenon, activity, ...</p> <p>Input: Batch of <code>{n_batch}</code> synset names with definitions.</p> <p>Output: JSON — <code>{`synset`: `ACCEPT/REJECT`, `meta-category`: ..., `explanation`: ...}</code></p>

Figure 20. Phase 0: Concept Preparation. (a) Supercategory assignment builds the taxonomy hierarchy. (b) Tangibility filter removes abstract/non-visual concepts from the candidate pool.

Phase 1: MC-VQA Task Generation (confusing-category variant shown; single-image & multi-image variants analogous)
<p>Role: Domain expert for <code>{category_name}</code>.</p> <p>Goal: Design challenging, domain-specific MC questions testing fine-grained visual knowledge of <code>{category_name}</code>.</p> <p>§1 Provided Information:</p> <ul style="list-style-type: none"> – Image of <code>{category_name}</code>; <code>{category_def}</code> – Textual list of professional visual attributes (§9) – List of visually confusing categories: <code>{confusing_categories_list}</code> <p>§2 Your Task:</p> <ol style="list-style-type: none"> 1. Analyze image, attribute list, and confusing categories (internal only). 2. Generate <code>{num_questions}</code> MC tasks; craft an internal plan for difficulty. 3. Web search for additional information about confusing categories as needed. <p>§3 Question Content: Visual focus, evidence-based, domain attributes only. <i>Task Definition:</i> <code>{task}</code> (one of 7 task types from Figure 15)</p> <p>§4 Choice Format: 5 choices (A–E); exactly 1 correct + 4 distractors.</p> <p>§5 Strict Constraints:</p> <ul style="list-style-type: none"> • <i>Visual Verification</i> — all options must reference visible/verifiable features; image takes priority over text. • <i>Professional Nomenclature</i> — formal domain-specific terms only; no colloquial hints. • <i>Content Restrictions</i> — avoid revealing <code>{category_name}</code> or broader classification; no quantitative measurements. • <i>Anti-Shortcut</i> — cannot be answered by elimination or text knowledge alone. <p>§6 Distractor Requirements: Visual similarity, terminology proximity, expert-level challenge, plausibility, category integration, task-specific challenges, knowledge–image conflict.</p> <p>§7 Output: JSON list — <code>[{`question`, `options`: [A...E], `answer`, `strategy`}]</code></p> <p>§8 Exception Handling: Return error JSON if image blurred/irrelevant, attributes corrupted, etc.</p> <p>§9 Attribute list for <code>{category_name}</code>: <code>{attribute_list}</code></p>

Figure 21. Phase 1: Main MC-VQA task generation prompt (confusing-category variant). The prompt injects a task-type definition (§3) from one of 7 task types. Single-image and multi-image variants share the same structure, differing only in §1 inputs and category references.

Phase 1 (Challenge): Iterative Difficulty Refinement (2 iterations; single/conf/MI variants)

Role: Expert biologist, evaluator and refiner of AI-generated visual assessments. *Has access to web search.*

Goal: Enhance difficulty of existing MC-VQA tasks so they are significantly more challenging for both expert humans and advanced MLLMs.

§1 Provided Information: Image of `{category_name}`, attribute list (§7), existing Q&A tasks (§8), web search, original task definition: `{task}`.

§2 Your Task:

1. Review image, attributes, and existing tasks; create internal mental description.
2. Plan to strategically increase difficulty following all requirements (§3–5).
3. Web search to connect with confusing categories, incorporate subtle visual info, ensure factual correctness.
4. Refine tasks: maximize difficulty by revising questions and options (especially distractors).

§3 Core Refinement Criteria:

- A. Maximize distractor plausibility & confusion.
- B. Apply task-specific strategies + broader misleading principles: `{challenge}`
- C. Prevent short-cuts: minimal redundancy, consistent format across options, require visual inspection.
- D. Maintain original task constraints and definition alignment.

§4 Constraints: (*same as generation: visual verification, nomenclature, content restrictions*)

Meta-Strategy Injected via `{challenge}`: *Subtle detail inaccuracies; attribute misattributions; plausible but absent features; incorrect spatial/relational descriptions; near-synonym substitutions; phylogenetic neighbor confusion; trait value shifts; contextual re-framing.*

Output: JSON with `refinement_notes` explaining changes (iteration 0: 4 choices; iteration 1: 5 choices).

Figure 22. Phase 1 (Challenge): Iterative difficulty refinement. Applied twice ($k=0, 1$) after initial generation. The meta-strategy list is injected via the `{challenge}` variable. Each iteration receives the output of the previous round as §8.