Do You See Me : A Multidimensional Benchmark for Evaluating Visual Perception in Multimodal LLMs

Anonymous ACL submission

Abstract

002

006

800

012

017

022

024

027

035

037

039

041

043

045

Multimodal Large Language Models (MLLMs) have demonstrated promising reasoning capabilities in diverse domains, yet their visual perception skills remain a critical bottleneck. In this study, we first investigate the impact of visual perception errors on visual reasoning questions by analyzing the performance of the model on 150 questions. Our findings reveal that incorrect answers often stem from failures in visual perception. In addition, some correct answers arise from hallucinated visual details. Towards this end, we introduce Do You See Me, a comprehensive benchmark dataset designed to evaluate visual perception capabilities in MLLMs. Drawing from established principles of human psychology, our dataset comprises 1,262 programmatically generated images that systematically assess different dimensions of visual understanding. Our benchmark consists of seven perception-focused subtasks, each designed with control parameters to modulate task complexity. Additionally, it can be easily extended for new perception tasks and varying complexities. We evaluate multiple state-of-the-art closed source and open source MLLMs and we conduct a comprehensive human study where participants are instructed to complete visual perception tasks and rate the difficulty of each test sample as easy, medium, or hard. Results indicate that MLLMs perform poorly on visual perception tasks, achieving less than 50% accuracy on most subtasks, whereas, humans perform with an average accuracy of around 95%. Furthermore, as task complexity increases, MLLM performance declines drastically, while human performance remains stable. In one of our subtasks, when analyzing performance across human-rated difficulty levels, we observe a progressive widening of gap between human and MLLM accuracy - starting at 12% for easy samples and expanding to 45% for those rated as difficult. This pattern of deteriorating MLLM performance relative to consistent human accuracy is seen across

all seven subtasks, highlighting a critical need	(
for enhanced visual perception capabilities in	(
MLLMs.	(

051

056

057

059

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

076

077

079

081

1 Introduction

Recent advances in Multimodal Large Language Models (MLLMs) have led to impressive performance in diverse areas, including mathematical problem-solving (Lu et al., 2023; Zhang et al., 2024a), physics-based questions (Saikh et al., 2022), and logical reasoning challenges (Xiao et al., 2024). Multiple visual question-answering (VQA) tasks have emerged as a critical benchmark for evaluating artificial intelligence systems, particularly in assessing their ability to combine visual perception with higher-level reasoning (Yue et al., 2023). While the long-term goal is to achieve superhuman reasoning across these tasks, an essential prerequisite is the ability to accurately perceive and interpret visual diagrams or figures-a capability rooted in fundamental spatial and form-recognition skills such as visual discrimination and spatial orientation.

Fig. 1 shows MLLM responses to a logical reasoning question based on pattern completion. We note that both Claude Sonnet-3.5 (Anthropic, 2024) and GPT-40 (OpenAI et al., 2024a) show visual perception errors in their reasoning steps for finding the correct answer with Claude Sonnet-3.5 even reaching the correct final answer despite improper visual grounding. Recent studies have highlighted incorrect or incomplete visual perception as a primary source of error when MLLMs attempt to answer visually grounded reasoning questions (Lu et al., 2023). This observation underscores the importance of thoroughly understanding the visualperception capabilities and failures even when a model's overall reasoning might appear correct. Visual reasoning itself can be viewed as a composition of three distinct processes: question comprehension (textual or visual), visual perception, and



Figure 1: Visual Misinterpretations in Popular Multimodal LLMs

reasoning. Because accurate perception is often a prerequisite for coherent reasoning, a model's success in high-level tasks might be misleading if it masks underlying weaknesses in visual skills.

086

087

090

100

102

103

104

106

108

109

110

111

112

113

114

115

116

117

118

Recent benchmarks such as MMVP (Tong et al., 2024b) and CV-Bench (Tong et al., 2024a) do provide insights into general visual perception abilities in MLLMs, but they depend heavily on detailed human annotations-making them difficult to scale into larger benchmarks. In contrast, tests for human visual perception (e.g., the Test of Visual Percption Skills (Gardner, 1988)) commonly employ structured difficulty levels. For instance, in shape discrimination tasks, the number and arrangement of shapes gradually increase in complexity, and participants must identify the presence of a specific shape. As the test becomes more challenging, performance eventually reaches a ceiling, revealing the upper limit of a subject's visual perception capabilities. A similar approach for MLLMs-where stimuli become progressively harder-would help pinpoint the exact complexity threshold beyond which these models fail to respond correctly, thus providing a more precise measure of their perceptual strengths and weaknesses. This understanding can be crucial in designing more targeted pretraining datasets, ensuring that models receive the necessary variety and difficulty of visual samples. It can also guide architectural refinements, ultimately leading to more robust and accurate multimodal reasoning systems.

Drawing on human psychology (Chalfant and Scheffelin, 1969), visual perception can be catego-

rized into five major interrelated skills: (1) visual discrimination, (2) visual figure-ground, (3) visual spatial, (4) visual closure, and (5) visual sequential memory. These categories provide a more nuanced framework to investigate where exactly models fail. Building on this perspective, our study addresses three primary research questions:

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

- **RQ1:** If a MLLM correctly answers a reasoning question based on a visual diagram, does it necessarily succeed on a closely related perception question about the same diagram?,
- **RQ2:** How effectively do current MLLMs perform on various dimensions of visual perception skills? and
- **RQ3:** What patterns emerge in the successes and failures of MLLMs across different visual perception tasks and their complexities, and how do these patterns compare to human performance?

To investigate **RQ1** (§2), we curate a joint visual perception-reasoning dataset comprising 150 unique images from logic-based Intelligence Quotient (IQ) tests. We opt for IQ-style questions because they focus on pattern recognition and "pure" reasoning with minimal domain-specific knowledge, making them well-suited for assessing raw reasoning abilities in MLLMs. Our dataset is constructed by merging samples from both Math-Vista (Lu et al., 2023) and LogicVista (Xiao et al., 2024) that feature logic-VQA tasks. For each image-question pair, we evaluate multiple MLLMs

on both perception-oriented and reasoning-oriented 150 questions. We then manually annotate each response, categorizing errors into four types: (1) visual perception error, (2) reasoning error, (3) arith-153 metic error, and (4) no error. Because a single 154 response can exhibit multiple error types, we as-155 sign all applicable labels to each response. This 156 comprehensive annotation process provides a more nuanced understanding of MLLM performance in 158 both perception and reasoning. 159

151

152

157

161

162

163

164

167

168

169

170

172

173

174

175

176

177

178

179

180

181

183

184

185

186

188

190

191

193

194

195

196

197

198

199

201

Next, to investigate RO2 and RO3 (§3), we introduce the Do You See Me benchmark, a programmatically generated test suite inspired by the five core dimensions of human visual perception. Our benchmark comprises 1,262 unique images designed to test visual perception, accompanied by 2,116 evaluation questions (see Appendix D for details). Our benchmark consists of seven subtasks designed to evaluate different aspects of visual perception. For every subtask, we define a set of control variables (see Appendix Table 4) that can be systematically adjusted, enabling us to increase or decrease task difficulty as needed. Because Do You See Me is fully programmatic, we can generate a large number of Image-QA pairs across varied difficulty levels. Our design approach facilitates a fine-grained analysis of MLLMs' perceptual capabilities, making it particularly valuable for pinpointing the specific visual skills at which MLLMs excel—and where they continue to struggle.

Our experiments on the joint perceptionreasoning dataset reveal that relying solely on finalanswer accuracy may obscure critical shortcomings. For example, Claude Sonnet-3.5, which was the best-performing model on reasoning questions with around 41% accuracy, still produced reasoning chains containing visual perception errors. Specifically, around 43% of the correctly answered samples exhibited visual perception errors. While, on the incorrectly answered samples, both reasoning and perception errors were found in equal number of samples.

On the **Do You See Me** benchmark, we observe that no single model outperforms all others across the various dimensions of visual perception. Further, we note that while humans perform well on all the seven subtasks with an average accuracy of 94.31%, the best performing MLLM exhibits an accuracy of 50.05%. Additionally, we observe that as we use increasingly difficult combination of control parameters, MLLMs show a drastic drop in performance (dropping to near zero accuracy), whereas human performance remains largely stable. Results on the proposed Do You See Me benchmark clearly indicate that MLLMs fare poorly on visual perception skills, and exhibit a large gap on performance when compared to humans. This suggests an urgent need to improve visual perception capabilities in MLLMs, independently from high-level reasoning.

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

237

238

239

240

241

242

243

244

245

247

248

249

250

In summary, our key contributions are:

- 1. A joint dataset where reasoning questions on logic-based VQA are paired with closely related visual perception questions.
- 2. An extensive error analysis of SOTA MLLMs for their performance on the joint dataset, including extensive manual annotation of reasoning chains to identify the errors.
- 3. A dynamic and synthetic test bench—Do You See Me—that systematically measures fundamental visual-perception skills inspired by established categories in human psychology.
- 4. A comprehensive evaluation of various closedsource and open-source MLLMs against human performance, uncovering patterns of success and failure, and underscoring the gap between "reasoning" and "seeing" in current MLLMs.

2 Preliminary Study

2.1 Joint Visual Perception and Reasoning Dataset

Most of the existing benchmarks assess MLLM's visual reasoning capabilities by solely relying on the final answer based accuracy. However, this approach can obscure the exact source of errors. In particular, three primary sources of error can arise: 1) Visual Perception- inaccuracies in identifying or interpreting elements in the provided image, 2) Reasoning- errors in the logical or conceptual steps used to arrive at the final answer, or 3) Arithmeticmistakes in performing numerical or algebraic calculations. To accurately distinguish between different error sources, it is essential to analyze not only final answers but also the reasoning chains. We introduce a joint perception-reasoning dataset specifically designed to separate visual perception errors from higher-level reasoning failures.

Why IQ-Type Questions? IQ-style diagrammatic questions primarily feature basic geometric shapes and patterns, minimizing reliance on domain-specific knowledge. This allows for a focused evaluation of visual perception and reasoning

- 251
- 25
- 25
- 255 256 257 258 259
- 260 261 262
- 266 267

264

268

269 270 271

273 274

275

27

278 279

- 28
- 28
- 282
- 20
- 284

2

- 2
- 2
- 289 290

skills without introducing extraneous complexity.

2.1.1 Data Collection

Our dataset is drawn from two established visual reasoning benchmarks: MathVista (Lu et al., 2023) and LogicVista (Xiao et al., 2024). We selected logic-based tasks centered around geometric shapes and pattern recognition from:

- *IQtest* subset of MathVista (focusing on spatial and pattern-based problems).
- *Diagrams* subset of LogicVista (pattern completion tasks).

These subsets feature universally understood shapes in controlled layouts, allowing systematic evaluation of perception and reasoning. We curated 15 problems from MathVista's *IQtest* and 135 from LogicVista's *Diagrams*, yielding a total of 150 examples in our final dataset. Table 1 summarizes the distribution.

Table 1: Distribution of examples across MathVista and LogicVista. "IQ/Logic Qs." refers to pattern-based or spatial reasoning questions.

Dataset	Original Size	IQ/Logic Qs	Selected
$MathVista_{mini}(IQtest)$	1000	37	15
LogicVista (Diagrams)	448	223	135

2.1.2 Data Annotation

We extend each original problem (I, R, A_R) where I is the image, R is the reasoning question, and A_R is the corresponding ground-truth answer—by adding a visual perception question Pwith its ground-truth A_P . The extended sample is thus: $(I, (R, A_R), (P, A_P))$. Five early-stage researchers proposed potential perception questions directly relevant to each reasoning question (e.g., "How many triangles are in the figure?"). Two senior annotators then refined these to ensure each perception question was unambiguous and closely tied to the underlying reasoning skill. Ten ambiguous queries were discarded and replaced, leading to a total of 150 finalized samples.

3 Do You See Me

Human psychology systematically categorizes human visual perception as a combination of five core abilities (Chalfant and Scheffelin, 1969):

• Visual Discrimination: Ability to recognize dominant features (e.g., position, shape, form, color).

Sub-Divisions of the Motor-Free Visual Perception Test						
		S Z Z Z Ú	ZZZZZZ	-		
Form Discrimination	Visual Figure Ground	Visual Closure	Visual Spatial	Visual-Sequential Memory		

Figure 2: A collection of samples for each subdivision of the TVPS test (Gardner, 1988)

• **Visual Figure-Ground:** Ability to distinguish the main object from its background.

291

292

293

294

295

297

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

- **Visual Memory:** Ability to remember sequences of presented images.
- **Visual Closure:** Ability to complete partially obscured shapes.
- **Visual Spatial:** Ability to perceive positions of objects relative to oneself and to other objects.

Assessments such as the Test of Visual Perception Skills (TVPS) (Gardner, 1988) and Motor-Free Visual Perception Test (MVPT) (Colarusso, 2003) systematically evaluate these abilities through structured questionnaires. Although MLLMs differ from human vision, these categories offer a solid framework for evaluating visual perception in MLLMs. Building on these principles, we introduce the Do You See Me benchmark, a fully automated test suite that dynamically generates both images and perception-focused questions (VPQA) of incremental difficulty. Our approach allows us to have a graded evaluation of MLLM performance over multiple perception dimensions (refer Appendix D for dataset distribution). Note that visual memory in addition to visual perception is also a test of short term visual memory in humans. Since the MLLMs currently have no memory, except via a textual description of the visual scene, we leave out visual memory from our benchmark. Our code is completely written in Python3¹, and uses SVG representation to generate the visual images. We open-source the synthetic data generation code, and dataset at 2 .

3.1 Visual Discrimination Tasks

We organize visual discrimination into four subtasks.

Shape Discrimination: This subtask assesses the ability to count specific shapes within a composite image. We use seven geometric shapes (*rectan*-

¹https://www.python.org/downloads/

²https://anonymous.4open.science/r/DoYouSeeMe-F52E/README.md



Figure 3: Do You See Me benchmark visual perception dimensions

329 gle, triangle, circle, pentagon, hexagon, octagon, star), each with a solid black border and transparent interior. Non-overlapping placements are 331 enforced using the Separating Axis Theorem (SAT) (Gottschalk et al., 1996). The difficulty is con-333 334 trolled by three parameters: the number of unique shapes (S), the maximum instances per shape (S_I) , and an overlap factor (α). Each generated image is paired with questions (for example, "How many circles are in the image?") and its programmatic 338 339 ground truth answer.

Joint Shape-Color Discrimination: This subtask focuses on compositional counting queries (e.g., "Count all red triangles"), using six shapes (*star, triangle, pentagon, hexagon, octagon, cross*) and eight colors (*red, green, blue, orange, purple, black, gray, yellow*). We use non-overlapping shapes in the task to avoid ambiguity in shape recognition. Difficulty is governed by the number of unique shapes (S) and the number of unique colors (C).

341

347

349Letter Discrimination: Letters are constructed350from a 5×7 LED-style binary matrix, testing a351model's ability to recognize block-based letters.352Task difficulty depends on the block spacing fac-353tor (β), the color contrast (ΔC) between the let-354ter blocks and the background, and the number of355letters (N) present. Higher block spacing (β) or356lower color contrast (ΔC) increases difficulty, as357does adding more letters in the scene.

Visual Form Constancy: This subtask checks if
a model can recognize a target pattern after shape
substitutions, rotations, and scaling. The target pat-

tern is constructed from simple primitives (*circle*, 361 square, line, triangle). We create three variants of 362 this pattern by applying a shape substitution factor (ssf), a scaling factor (α) , and a rotation factor 364 (θ_r) . The MLLM must identify the option that 365 exactly matches the target's arrangement. 366

367

368

370

371

373

374

375

376

378

379

381

382

383

385

3.2 Visual Spatial

This task uses one or more $H \times W$ grids with each cell containing one of three shapes (*circle, square, triangle*) rendered in solid black or outlined form. The parameters are the grid dimension (D) and the number of grids (G). The visual query is contructed instructing the MLLM to locate shapes by row-column coordinates or counting shapes relative to a reference position in the grid (e.g., "How many solid circles are above the triangle in row 3, column 2?").

3.3 Visual Figure-Ground

We extend the visual form constancy framework by introducing distracting background elements. Two parameters control the task complexity: 1) number of shapes (N) and, the background distraction factor (bdf). The bdf determines the amount of noise added to the background. The task of the MLLM is to distinguish the target pattern from other candidates under presence of background noise.

3.4 Visual Closure

390

391

396

400

401

402

403 404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428 429

430

431

432

433

This subtask evaluates whether the MLLM can match an incomplete shape with the corresponding fully formed target. We define seven basic shapes (*capsule, star, hexagon, circle, pentagon, rectangle, triangle*), remove some edges to create the incomplete target, and then produce three "noisy" distractors by randomly distorting vertex positions. The model must select which incomplete shape correctly closes to the original target.

3.5 Synthetic Data Generation

We specify a range of control parameters for each subtask (see Table 4) and generate all possible parameter combinations. Each combination yields ten images, each paired with multiple questions. The parameter limits are identified by conducting a preliminary study over visual perception accuracy for GPT-40 (OpenAI et al., 2024b). Design of each subtask is elaborated in Appendix B.

3.6 Human Performance Benchmarking

We compare modern MLLMs with seven human participants on the **Do You See Me** dataset. For each subtask, subjects answered two randomly selected questions per parameter combination and rated difficulty. A brief calibration phase preceded each subtask to reduce bias. Human accuracy, computed by matching typed answers to ground truth, provides a baseline for model comparisons (refer Appendix E for more details).

4 Experiments

4.1 Evaluation Protocol

Recent Large Language Models (LLMs) and Multimodal Large Language Models (MLLMs) are increasingly instructed to produce extended textual outputs rather than concise responses, making earlier rule-based or template-matching methods (Lu et al., 2022) difficult to apply. Inspired by recent benchmarks for MLLMs (Lu et al., 2023; Zhang et al., 2024a), we employ an expert LLM to evaluate answers. Our framework proceeds in three stages. In the first stage, a MLLM generates a detailed response according to a predefined template (see Appendix G), which includes the task description, the question, and possible choices. Next, an answer extractor (Appendix G), based on GPT-40 (OpenAI et al., 2024b), parses these extended outputs to yield a concise answer. Prior work has



Figure 4: Radar chart illustrating MLLM performance on the seven subtasks in **Do You See Me** benchmark dataset.

shown that such an expert LLM can extract the correct answer with near 100% accuracy (Lu et al., 2023). Finally, the extracted text is standardized (e.g., reduced to multiple-choice labels or numeric values), and performance metrics are computed. Since the **Do You See Me** dataset contains both multiple-choice (textual) and free-response (numeric) questions, accuracy is used as a measure of performance.

434

435

436

437

438

439

440

441

442

443

4.2 Experimental Setup

We evaluate our Do You See Me dataset using 444 both closed-source and open-source Multimodal 445 Large Language Models (MLLMs). For closed-446 source models, we use GEMINI-1.5 FLASH (Team 447 et al., 2024), GPT-40 (OpenAI et al., 2024a), and 448 CLAUDE-SONNET-3.5 (Anthropic, 2024). For 449 open-source models, our assessment includes 450 LLAMA-3.2-11B-VISION (Grattafiori et al., 2024), 451 LLAMA-3.2-90B-VISION (Grattafiori et al., 2024), 452 DEEPSEEK-VL2-SMALL-3B (Wu et al., 2024b), 453 DEEPSEEK-VL2-TINY-1B (Wu et al., 2024b), 454 QWEN2.5-VL-7B-INSTRUCT (Wang et al., 2024), 455 and INTERNVL2.5-8B (Chen et al., 2025). Each 456 model is given identical inputs, both the same vi-457 sual content and uniform textual prompts, to ensure 458 fair comparison. The closed-source MLLMs are 459 accessed through their respective proprietary APIs. 460 In contrast, the open-source MLLMs are run locally 461 with consistent hyperparameter settings, including 462 batch size and decoding parameters. 463

4.3 Experimental Results

464

465

466

467

468

469

470

471

472

473

474

475

476 477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

502

4.3.1 Joint Perception-Reasoning Dataset

In this section, we evaluate MLLMs on a **joint perception-reasoning dataset**, where each sample contains an image (I), a perception question ((P)), and a reasoning question (R). Each model is separately asked to answer both questions with explicit instructions tp respond with a detailed chainof-thought. We use an LLM-based grader (see Section 4.1) to score the answers on correctness.



Figure 5: *MLLM performance on joint perceptionreasoning questions*

Perception vs. Reasoning Performance. Figure 5 illustrates the distribution of joint performance on the perception and reasoning questions for each image. The stacked bars show four categories: 1) both perception and reasoning correct, 2) both incorrect, 3) correct reasoning but incorrect perception, and 4) correct perception but incorrect reasoning. Claude Sonnet-3.5 has the largest fraction of samples with both perception and reasoning correct, whereas DeepSeek-VL2-Small-3B has the highest rate of overall failure. Interestingly, there are some instances where a model provides a correct reasoning answer despite failing the corresponding perception question. This can be caused by one of three things: 1) the perception question being only loosely related to the reasoning question, 2) the reasoning question resulting in better interpretation of the figure than a direct visual perception question leading to a correct final answer, or 3) the model arriving at the correct reasoning conclusion through hallucinatory or shortcut strategies, not genuine understanding. Thus, a closer, fine-grained analysis is needed to pinpoint which factor dominates.

Error Breakdown in Reasoning Chains. To gain deeper insights, we manually annotated answers from top-2 models (Claude Sonnet-3.5, GPT-40), labeling each error as *visual perception, reasoning*, or *arithmetic* (as defined in Section 2). Fig-

ure 6 presents the error distribution for Claude Sonnet-3.5 (see Appendix C for GPT4o results). We observe that, in the case where the final reasoning answer was determined to be incorrect, both visual perception and reasoning errors were a cause. Whereas, on cases where the final answer was correct, the detailed model response contained around 42.9% visual perception errors, with only 14.3% samples were answered error free. Clearly, these findings underscore the importance of improving visual perception in MLLMs. They also reinforce that success on a final question does not necessarily imply robust grounding in visual details or error free reasoning. 503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

529

530

531

532

533

534

535

536

537



(a) Error distribution for correct final answers.



(b) Error distribution for incorrect final answers.

Figure 6: *Error distributions in reasoning chain for* visual reasoning questions (Claude Sonnet-3.5).

4.3.2 Do You See Me

In this section, we present a comparative evaluation of multiple MLLMs (and human experts) on a set of visual reasoning tasks from the Do You See Me benchmark. These tasks assess a broad range of perceptual and interpretive skills, including visual closure, form constancy, figure-ground separation, spatial reasoning, color/shape disambiguation, and letter disambiguation. Overall Performance Across Dimensions Figure 4 compares the accuracy of all models across the various subtasks. Each radial spoke corresponds to a distinct visual skill (e.g., visual closure or visual spatial), while the distance from the center (0%-100% accuracy) denotes performance. Consistent with the quantitative results in Table 2, we observe that current MLLMs still lag behind human performance by a large margin across all subtasks. Despite this gap, several interesting observations emerge: Most MLLMs are relatively strong visual form constancy. For instance, Claude Sonnet-3.5 achieves an accu-



Figure 7: *MLLMs vs Human on Increasingly Difficult Combinations of Control Parameters for Visual Form Constancy Task.*

racy of 91.8%, the highest on form constancy. Further, MLLMs also perform relatively well on joint color-shape disambiguation with Gemini-1.5 Flash and Qwen2.5-VL-7B-Instruct performing the best. Finally, we observe that no single model is the best across all dimensions. Different tasks highlight distinct top performers; for example, Claude Sonnet-3.5 leads on form constancy, whereas Gemini-1.5 Flash and Qwen2.5-VL-7B-Instruct excel at color-shape tasks. Notably, the radar chart also reveals stark performance drops when tasks involve noise, occlusion, or overlap. For example, in visual figure-ground addition of noise in comparison to the form constancy task results in a gap of roughly 50% between top-performers. A similar trend appears between the shape discrimination task (which involves overlapping black outlines) and the joint shape-color task (no overlap). Thus, it is plausible that MLLM accuracy deteriorates significantly as tasks incorporate occlusion or noise around critical visual areas.

538

539

540

542

543

544

547

549

551

554

555

557

558

559

563

564

567



Figure 8: Human vs MLLM Performance on Human Rated Difficulty Levels

Model Performance Over Increasing Difficulty In §3.6, we described our human study on the **Do You See Me** dataset, where participants provided *perceived difficulty* ratings for each combination of control parameters across all subtasks. We used majority voting to classify each combination as *easy*, *medium*, or *hard*. Next, we partitioned the accuracy results of each MLLM according to these levels. For clearer comparison, we grouped MLLMs into two categories-closed-source and open-source-and then averaged their performance within each difficulty level, alongside human performance. Figure 8 shows the results for two example tasks: visual form constancy and letter disambiguation. On visual form constancy-where MLLMs generally perform better than on other subtasks-both open- and closed-source models show progressively lower accuracies as human-rated task difficulty increases. Notably, the performance gap compared to humans widens: for closed-source models, it grows from about 12% at Easy to 20% at Medium and 45% at Hard. Closed-source models also outperform open-source models across all difficulty levels. Meanwhile, on the letter disambiguation task, as soon as the difficulty reaches Medium, both closed- and open-source models effectively fail (performance is near zero), whereas human accuracy remains high. We observed similar trends in the other subtasks as well (see Appendix C).

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

5 Conclusion

In this study, we examined the visual perception capabilities of MLLMs and their impact on logical reasoning. Our preliminary analysis revealed that many reasoning errors in MLLMs originate from incorrect visual perception, while some correct answers result from hallucinated visual details. These findings highlight the critical role of perception in multimodal reasoning and motivate the need for a systematic evaluation of MLLM visual skills. Motivated by these findings, we introduced the Do You See Me —a programmatic and scalable benchmark with a collection of seven visual perception subtasks inspired by human psychology.

Experimental results demonstrated that MLLMs struggle considerably on these subtasks compared to humans, particularly when complexity of the subtasks increase. While human performance remained robust across varying difficulty levels, MLLMs showed rapid declines in accuracy, reinforcing the need for improved visual processing. Overall, our study highlights the importance of developing more perceptually grounded MLLMs to reduce hallucinations and ensure reliable performance on visual reasoning and related tasks. Our approach of programmatic, scalable and complexity controlled data creation is not only suitable for evaluations, but for synthetic training data creation as well.

6 Limitations

617

Our work has a few limitations that we acknowl-618 edge and plan to address in future research. First, 619 the size of our joint perception-reasoning dataset is relatively small. However, we have made every effort to include all possible samples where it was feasible to generate non-ambiguous and correlated visual perception questions. To address this limitation, we plan to employ LLM + 625 Image-Diffusion techniques in the future to generate a more diverse and controlled format of perception+reasoning questions, thus expanding our dataset. Second, while our dataset considers nonreal-world settings for evaluating visual perception, we found this approach extremely useful for generating a vast number of diverse examples with a broad difficulty range. It is worth noting that most human visual perception tests follow a similar format. In future work, we aim to explore the use of 3D tools such as Unity3D $^{\rm 3}$ and Blender $^{\rm 4}$ to generate realistic 3D questions based on the same 637 principles presented in our dataset, thus enhancing the ecological validity of our experiments. Third, in the current setup, we have restricted our visual perception prompts to the English language only, including the letter disambiguation task. This deci-643 sion was made in the interest of managing the overall cost of benchmarking closed-source MLLMs. 644 However, we recognize the importance of language diversity and plan to expand our coverage to other non-English languages in future iterations of our 647 work. Overall, we believe that our work provides a valuable contribution to the understanding of MLLM capabilities in visual perception tasks and lays the foundation for future research in this area.

References

654

655

656

657

663

- Anthropic. 2024. Claude 3.5 sonnet. https://www. anthropic.com. Large language model, released October 2024.
- Ted Brown and Lisa Peres. 2018. An overview and critique of the test of visual perception skills – fourth edition (tvps-4). *Hong Kong Journal of Occupational Therapy : HKJOT*, 31:59 – 68.
- James C Chalfant and Margaret A Scheffelin. 1969. Central processing dysfunctions in children: A review of research.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye,

Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Botian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu, Huipeng Deng, Jiaye Ge, Kai Chen, Kaipeng Zhang, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. 2025. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *Preprint*, arXiv:2412.05271. 665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

- Ronald P. Colarusso. 2003. Mvpt-3: Motor-free visual perception test.
- Morrison F. Gardner. 1988. Tvps, test of visualperceptual skills (non-motor).
- Stefan Gottschalk, Ming C Lin, and Dinesh Manocha. 1996. Obbtree: a hierarchical structure for rapid interference detection. *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew

³https://unity.com/products/unity-engine

⁴https://docs.blender.org/manual

726 Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, 727 Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj 734 Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, 736 Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-737 nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-740 hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-741 hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye 742 743 Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten 744 Sootla, Stephane Collot, Suchin Gururangan, Syd-745 746 ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias 747 Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-751 ginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-753 ney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-754 755 feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing 758 Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, 761 Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-770 dan, Beau James, Ben Maurer, Benjamin Leonhardi, 771 Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-772 773 cock, Bram Wasti, Brandon Spence, Brani Stojkovic, 774 Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, 775 Changkyu Kim, Chao Zhou, Chester Hu, Ching-776 777 Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-778 ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, 779 Daniel Kreymer, Daniel Li, David Adkins, David 780 Xu, Davide Testuggine, Delia David, Devi Parikh, 781 Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, 783 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank 787 788 Seide, Gabriela Medina Florez, Gabriella Schwarz, 789 Gada Badeer, Georgia Swee, Gil Halpern, Grant

Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,

790

791

792

793

794

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

- 855 856
- 85
- 80 86
- 861 862
- 863 864 865
- 866 867 868
- 8
- 871 872
- 8
- 874 875
- 877 878 879

81 81 81

8 8 8

900

904

905

906

910

911

901 902 903

907 908 909

- Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.
- D.D. Hammill, J.K. Voress, and N.A. Pearson. 2016. *DTVP-3*. Manual Moderno.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. 2022. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International Conference on Machine Learning*.
- Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. In *Annual Meeting of the Association for Computational Linguistics*.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chun yue Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations*.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and A. Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *ArXiv*, abs/2209.09513.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson,

Chak Li, Chan Jun Shern, Channing Conger, Char-912 lotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, 913 Chong Zhang, Chris Beaumont, Chris Hallacy, Chris 914 Koch, Christian Gibson, Christina Kim, Christine 915 Choi, Christine McLeavey, Christopher Hesse, Clau-916 dia Fischer, Clemens Winter, Coley Czarnecki, Colin 917 Jarvis, Colin Wei, Constantin Koumouzelis, Dane 918 Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, 919 David Carr, David Farhi, David Mely, David Robin-920 son, David Sasaki, Denny Jin, Dev Valladares, Dim-921 itris Tsipras, Doug Li, Duc Phong Nguyen, Duncan 922 Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-923 dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, 924 Eric Kramer, Eric Peterson, Eric Sigler, Eric Wal-925 lace, Eugene Brevdo, Evan Mays, Farzad Khorasani, 926 Felipe Petroski Such, Filippo Raso, Francis Zhang, 927 Fred von Lohmann, Freddie Sulit, Gabriel Goh, 928 Gene Oden, Geoff Salmon, Giulio Starace, Greg 929 Brockman, Hadi Salman, Haiming Bao, Haitang 930 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, 931 Heather Whitney, Heewoo Jun, Hendrik Kirchner, 932 Henrique Ponde de Oliveira Pinto, Hongyu Ren, 933 Huiwen Chang, Hyung Won Chung, Ian Kivlichan, 934 Ian O'Connell, Ian O'Connell, Ian Osband, Ian Sil-935 ber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya 936 Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, 937 Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub 938 Pachocki, James Aung, James Betker, James Crooks, 939 James Lennon, Jamie Kiros, Jan Leike, Jane Park, 940 Jason Kwon, Jason Phang, Jason Teplitz, Jason 941 Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-942 avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui 943 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, 944 Joaquin Quinonero Candela, Joe Beutler, Joe Lan-945 ders, Joel Parish, Johannes Heidecke, John Schul-946 man, Jonathan Lachman, Jonathan McKay, Jonathan 947 Uesato, Jonathan Ward, Jong Wook Kim, Joost 948 Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, 949 Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, 950 Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai 951 Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin 952 Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, 953 Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, 954 Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle 955 Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-956 ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia 957 Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-958 ian Weng, Lindsay McCallum, Lindsey Held, Long 959 Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kon-960 draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, 961 Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine 962 Boyd, Madeleine Thompson, Marat Dukhan, Mark 963 Chen, Mark Gray, Mark Hudnall, Marvin Zhang, 964 Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, 965 Max Johnson, Maya Shetty, Mayank Gupta, Meghan 966 Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao 967 Zhong, Mia Glaese, Mianna Chen, Michael Jan-968 ner, Michael Lampe, Michael Petrov, Michael Wu, 969 Michele Wang, Michelle Fradin, Michelle Pokrass, 970 Miguel Castro, Miguel Oom Temudo de Castro, 971 Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-972 nal Khan, Mira Murati, Mo Bavarian, Molly Lin, 973 Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-974 talie Cone, Natalie Staudacher, Natalie Summers, 975

Natan LaFontaine, Neil Chowdhury, Nick Ryder, 976 Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, 977 Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel 978 Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra-987 jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-991 dani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu 997 Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao 1001 Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunninghman, Thomas 1003 Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor 1004 Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, 1006 Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Oi, Veit Moeller, 1007 Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne 1008 Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, 1009 1010 Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, 1011 Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen 1012 He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and 1013 Yury Malkov. 2024a. Gpt-4o system card. Preprint, arXiv:2410.21276. 1014

> OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott

1015

1016

1017 1018

1019

1020

1021

1022 1023

1024

1025

1026

1027 1028

1029 1030

1031

1032

1033

1034

1035

1036

1037

1038

Gray, Ryan Greene, Joshua Gross, Shixiang Shane 1039 Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, 1040 Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, 1042 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin 1043 Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-1046 woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-1047 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, 1048 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, 1049 Christina Kim, Yongjik Kim, Jan Hendrik Kirch-1050 ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, 1051 Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-1052 stantinidis, Kyle Kosic, Gretchen Krueger, Vishal 1053 Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan 1054 Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, 1057 Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie 1059 Mayer, Andrew Mayne, Bob McGrew, Scott Mayer 1060 McKinney, Christine McLeavey, Paul McMillan, 1061 Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 1064 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, 1066 Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, 1067 Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambat-1069 tista Parascandolo, Joel Parish, Emy Parparita, Alex 1070 Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-1071 man, Filipe de Avila Belbute Peres, Michael Petrov, 1072 Henrique Ponde de Oliveira Pinto, Michael, Poko-1073 rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-1074 ell, Alethea Power, Boris Power, Elizabeth Proehl, 1075 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, 1076 Cameron Raymond, Francis Real, Kendra Rimbach, 1077 Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-1078 der, Mario Saltarelli, Ted Sanders, Shibani Santurkar, 1079 Girish Sastry, Heather Schmidt, David Schnurr, John 1080 Schulman, Daniel Selsam, Kyla Sheppard, Toki 1081 Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, 1083 Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin 1084 Sokolowsky, Yang Song, Natalie Staudacher, Fe-1085 lipe Petroski Such, Natalie Summers, Ilya Sutskever, 1086 Jie Tang, Nikolas Tezak, Madeleine B. Thompson, 1087 Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, 1088 Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-1089 lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, 1090 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, 1091 Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, 1092 CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-1093 ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, 1094 Clemens Winter, Samuel Wolrich, Hannah Wong, 1095 Lauren Workman, Sherwin Wu, Jeff Wu, Michael 1096 Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, 1097 Qiming Yuan, Wojciech Zaremba, Rowan Zellers, 1098 Chong Zhang, Marvin Zhang, Shengjia Zhao, Tian-1099 hao Zheng, Juntang Zhuang, William Zhuk, and Bar-1100 ret Zoph. 2024b. Gpt-4 technical report. Preprint, 1101 arXiv:2303.08774. 1102

- 1103 1104
- 1105 1106
- 1107 1108
- 1109 1110
- 1111 1112
- 1113
- 1114 1115
- 1116 1117

1118

1123

1124

- 1125 1126 1127 1128 1129 1130 1131 1132 1133
- 1134 1135 1136 1137 1138 1139 1140 1141

1142

1143 1144

1145

1146

1147

1148

1149

1150

1151

1152 1153

1154 1155

1156

1157

1158

1159

1160

1161

1162

1163

- Wujian Peng, Sicheng Xie, Zuyao You, Shiyi Lan, and Zuxuan Wu. 2023. Synthesize, diagnose, and optimize: Towards fine-grained vision-language understanding. 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 13279– 13288.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*.
 - Tanik Saikh, Tirthankar Ghosal, Amish Mittal, Asif Ekbal, and Pushpak Bhattacharyya. 2022. Scienceqa: a novel resource for question answering on scholarly articles. *International Journal on Digital Libraries*, 23:289 – 301.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, Andrea Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love, Paul Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, Tao Zhu, Kornraphop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, Manaal Faruqui, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal, Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vodrahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Siddhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo Xu, Raphael Lopez Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom, Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He, Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek, Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen,

Kaushik Shivakumar, Rishabh Agarwal, Harshal 1164 Godhia, Ravi Rajwar, Javier Snaider, Xerxes Doti-1165 walla, Yuan Liu, Aditya Barua, Victor Ungureanu, 1166 Yuan Zhang, Bat-Orgil Batsaikhan, Mateo Wirth, 1167 James Qin, Ivo Danihelka, Tulsee Doshi, Martin 1168 Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Ar-1169 jun Kar, Madhu Gurumurthy, Cheng Li, Ruoxin 1170 Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz, 1171 Nathan Lintz, Harsh Mehta, Heidi Howard, Mal-1172 colm Reynolds, Lora Aroyo, Quan Wang, Lorenzo 1173 Blanco, Albin Cassirer, Jordan Griffith, Dipanjan 1174 Das, Stephan Lee, Jakub Sygnowski, Zach Fisher, 1175 James Besley, Richard Powell, Zafarali Ahmed, Do-1176 minik Paulus, David Reitter, Zalan Borsos, Rishabh 1177 Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vi-1178 han Jain, Nikhil Sethi, Megha Goel, Takaki Makino, 1179 Rhys May, Zhen Yang, Johan Schalkwyk, Christina 1180 Butterfield, Anja Hauth, Alex Goldin, Will Hawkins, 1181 Evan Senter, Sergey Brin, Oliver Woodman, Mar-1182 vin Ritter, Eric Noland, Minh Giang, Vijay Bolina, 1183 Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid, 1184 Obaid Sarvana, David Silver, Alexander Chen, Lily 1185 Wang, Loren Maggiore, Oscar Chang, Nithya At-1186 taluri, Gregory Thornton, Chung-Cheng Chiu, Os-1187 kar Bunyan, Nir Levine, Timothy Chung, Evgenii 1188 Eltyshev, Xiance Si, Timothy Lillicrap, Demetra 1189 Brady, Vaibhav Aggarwal, Boxi Wu, Yuanzhong Xu, 1190 Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, 1191 Erica Moreira, Wojciech Stokowiec, Ross Hems-1192 ley, Dong Li, Alex Tudor, Pranav Shyam, Elahe 1193 Rahimtoroghi, Salem Haykal, Pablo Sprechmann, 1194 Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, 1195 Kalpesh Krishna, Xiao Wu, Alexandre Frechette, 1196 Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, 1197 Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao 1198 Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, 1199 HyunJeong Choe, Alex Tomala, Chalence Safranek-1200 Shrader, Nora Kassner, Mantas Pajarskas, Matt 1201 Harvey, Sean Sechrist, Meire Fortunato, Christina 1202 Lyu, Gamaleldin Elsayed, Chenkai Kuang, James 1203 Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Pe-1204 ter Humphreys, Kate Baumli, Connie Tao, Rajku-1205 mar Samuel, Cicero Nogueira dos Santos, Anders 1206 Andreassen, Nemanja Rakićević, Dominik Grewe, 1207 Aviral Kumar, Stephanie Winkler, Jonathan Caton, 1208 Andrew Brock, Sid Dalmia, Hannah Sheahan, Iain 1209 Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Fer-1210 yal Behbahani, Flavien Prost, Yanhua Sun, Artiom 1211 Myaskovsky, Thanumalayan Sankaranarayana Pillai, 1212 Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, 1213 Fabio Pardo, Xiaowei Li, Dan Horgan, Joe Stanton, 1214 Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu 1215 Wang, Basil Mustafa, Albert Webson, Hyo Lee, Ro-1216 han Anil, Martin Wicke, Timothy Dozat, Abhishek 1217 Sinha, Enrique Piqueras, Elahe Dabir, Shyam Upad-1218 hyay, Anudhyan Boral, Lisa Anne Hendricks, Corey 1219 Fry, Josip Djolonga, Yi Su, Jake Walker, Jane La-1220 banowski, Ronny Huang, Vedant Misra, Jeremy 1221 Chen, RJ Skerry-Ryan, Avi Singh, Shruti Rijh-1222 wani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, 1223 Romina Datta, Sumit Bagri, Arnar Mar Hrafnkels-1224 son, Marcello Maggioni, Daniel Zheng, Yury Sul-1225 sky, Shaobo Hou, Tom Le Paine, Antoine Yang, 1226 Jason Riesa, Dominika Rogozinska, Dror Marcus, 1227

1228 Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen 1229 Miller, Jeremy Greer, Lars Lowe Sjos, Azade Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu 1230 Jiang, Charlie Chen, Ruibo Liu, Tara Sainath, Maxim 1231 1232 Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh 1233 Newlan, Zeyncep Cankara, Soo Kwak, Yunhan Xu, 1234 Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihlas, Ada Ma, Juraj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, 1236 Sanjay Ganapathy, Ishita Dasgupta, Steph Hughes-1237 Fitt, Heng Chen, David Reid, Keran Rong, Hongmin 1238 1239 Fan, Joost van Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohananey, 1240 Anastasija Ilic, Taylor Tobin, John Wieting, Anna 1241 Bortsova, Phoebe Thacker, Emma Wang, Emily 1242 1243 Caveness, Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty, 1244 Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun 1246 Dai, Wenhao Jia, Matthew Wiethoff, Elnaz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel 1247 Gao, Golan Pundak, Susan Zhang, Michael Azzam, 1248 1249 Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma, Andy Swing, YaGuang Li, Chenxi 1250 Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary 1251 Nado, Ankesh Anand, Josh Lipschultz, Abhijit Kar-1252 1253 markar, Lev Proleev, Abe Ittycheriah, Soheil Hassas Yeganeh, George Polovets, Aleksandra Faust, 1254 1255 Jiao Sun, Alban Rrustemi, Pen Li, Rakesh Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh 1256 Baddepudi, Sebastian Krause, Emilio Parisotto, Radu 1257 1258 Soricut, Zheng Xu, Dawn Bloxwich, Melvin John-1259 son, Behnam Neyshabur, Justin Mao-Jones, Ren-1260 shen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal, Duc Dung Nguyen, James 1261 1262 Svensson, Le Hou, Sarah York, Kieran Milan, So-1263 phie Bridgers, Wiktor Gworek, Marco Tagliasacchi, 1264 James Lee-Thorp, Michael Chang, Alexey Guseynov, 1265 Ale Jakse Hartman, Michael Kwong, Ruizhe Zhao, 1266 Sheleem Kashem, Elizabeth Cole, Antoine Miech, Richard Tanburn, Mary Phuong, Filip Pavetic, Se-1267 1268 bastien Cevey, Ramona Comanescu, Richard Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, 1269 1270 Mariko Iinuma, Clara Huiyi Hu, Aurko Roy, Shaan 1271 Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel 1272 Saputro, Anita Gergely, Steven Zheng, Dawei Jia, 1273 Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina Samangooei, Mina 1274 Khan, Tomas Kocisky, Angelos Filos, Chintu Ku-1275 1276 mar, Colton Bishop, Adams Yu, Sarah Hodkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, 1277 1278 Yong Cheng, Adam Bloniarz, Jaehoon Lee, Pedram 1279 Pejman, Paul Michel, Stephen Spencer, Vladimir 1280 Feinberg, Xuehan Xiong, Nikolay Savinov, Char-1281 lotte Smith, Siamak Shakeri, Dustin Tran, Mary 1282 Chesus, Bernd Bohnet, George Tucker, Tamara von 1283 Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, 1284 Ambrose Slone, Kedar Soparkar, Disha Shrivastava, 1285 James Cobon-Kerr, Michael Sharman, Jay Pavagadhi, 1286 Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qiujia Li, Anton 1287 1288 Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao Tang, Eli 1289 1290 Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria 1291 Recasens, Guangda Lai, Alberto Magni, Nicola De

Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, 1292 Mostafa Dehghani, Jenny Brennan, Yifan He, Kelvin 1293 Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi 1294 Wu, Seb Arnold, Solomon Chang, Julian Schrit-1295 twieser, Elena Buchatskaya, Soroush Radpour, Mar-1296 tin Polacek, Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah Cogan, Aliaksei Severyn, Mohammad Saleh, 1299 Shantanu Thakoor, Laurent Shefey, Siyuan Qiao, 1300 Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao 1301 Zhang, Benjamin Lee, Paul Kishan Rubenstein, Gan 1302 Song, Tom Kwiatkowski, Anna Koop, Ajay Kan-1303 nan, David Kao, Parker Schuh, Axel Stjerngren, Gol-1304 naz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Fe-1305 lipe Tiengo Ferreira, Aishwarya Kamath, Ted Kli-1306 menko, Ken Franko, Kefan Xiao, Indro Bhattacharya, 1307 Miteyan Patel, Rui Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica Landon, Mara Finkelstein, Priya 1310 Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, 1311 Dennis Daun, Khuslen Baatarsukh, Matthew Tung, 1312 Wael Farhan, Henryk Michalewski, Fabio Viola, Fe-1313 lix de Chaumont Quitry, Charline Le Lan, Tom Hud-1314 son, Qingze Wang, Felix Fischer, Ivy Zheng, Elspeth 1315 White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwanicki, Michael Isard, 1317 Anna Bulanova, Lukas Zilka, Ethan Dyer, Deven-1318 dra Sachan, Srivatsan Srinivasan, Hannah Mucken-1319 hirn, Honglong Cai, Amol Mandhane, Mukarram 1320 Tariq, Jack W. Rae, Gary Wang, Kareem Ayoub, 1321 Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris 1322 Alberti, Dan Garrette, Kashyap Krishnakumar, Mai 1323 Gimenez, Anselm Levskaya, Daniel Sohn, Josip 1324 Matak, Inaki Iturrate, Michael B. Chang, Jackie Xi-1325 ang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian 1326 Hutter, Vahab Mirrokni, Nanxin Chen, Kaisheng 1327 Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty, 1328 Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat, 1329 Jasmine Liu, David Tao, Chloe Thornton, Tim Green, 1330 Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan 1331 Tan, Noah Fiedel, Hongkun Yu, Ed Chi, Alexan-1332 der Neitz, Jens Heitkaemper, Anu Sinha, Denny 1333 Zhou, Yi Sun, Charbel Kaed, Brice Hulse, Swa-1334 roop Mishra, Maria Georgaki, Sneha Kudugunta, 1335 Clement Farabet, Izhak Shafran, Daniel Vlasic, An-1336 ton Tsitsulin, Rajagopal Ananthanarayanan, Alen 1337 Carin, Guolong Su, Pei Sun, Shashank V, Gabriel 1338 Carvajal, Josef Broder, Iulia Comsa, Alena Repina, 1339 William Wong, Warren Weilun Chen, Peter Hawkins, 1340 Egor Filonov, Lucia Loher, Christoph Hirnschall, 1341 Weiyi Wang, Jingchen Ye, Andrea Burns, Hardie 1342 Cate, Diana Gage Wright, Federico Piccinini, Lei 1343 Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizh-1344 skaya, Ashwin Sreevatsa, Shuang Song, Luis C. 1345 Cobo, Anand Iyer, Chetan Tekur, Guillermo Gar-1346 rido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven 1347 Zheng, Hui Li, Ananth Agarwal, Christel Ngani, 1348 Kati Goshvadi, Rebeca Santamaria-Fernandez, Woj-1349 ciech Fica, Xinyun Chen, Chris Gorgolewski, Sean 1350 Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, 1351 Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian 1352 Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan 1353 Yuan, Florian Luisier, Alexandra Chronopoulou, Sal-1354 vatore Scellato, Praveen Srinivasan, Minmin Chen, 1355

Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, Achintya Singhal, Thang Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, Johnson Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li, Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Kopparapu, Francoise Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway, Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu, Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Milad Nasr, Ilia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy, Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Villela, Haibin Zhang, Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mahmoud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex Grills, Joseph Pagadora, Tsendsuren Munkhdalai, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihajlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto, Solomon Kim, William Zeng, Ken Durden, Priya Ponnapalli, Tiberiu Sosea, Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam, Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Bartek Perz, Wooyeol Kim, Nandita Dukkipati, Anthony Baryshnikov, Christos Kapla-

1356

1357

1358

1360

1364

1365

1366 1367

1368

1371

1374

1375

1376

1378

1379

1381

1383

1384

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

nis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, 1420 Diego de Las Casas, Harry Askham, Kathryn Tun-1421 yasuvunakool, Felix Gimeno, Siim Poder, Chester 1422 Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek 1423 Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai, 1424 Toby Shevlane, Christina Kouridi, Drew Garmon, 1425 Adrian Goedeckemeyer, Adam R. Brown, Anitha Vi-1426 jayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, 1427 Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep 1428 Kumar, Wei Chen, Courtney Biles, Garrett Bingham, 1429 Evan Rosen, Lisa Wang, Qijun Tan, David Engel, 1430 Francesco Pongetti, Dario de Cesare, Dongseong 1431 Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, 1432 Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aha-1433 roni, Trieu Trinh, Jessica Lo, Norman Casagrande, 1434 Roopali Vij, Loic Matthey, Bramandia Ramadhana, 1435 Austin Matthews, CJ Carey, Matthew Johnson, Kre-1436 mena Goranova, Rohin Shah, Shereen Ashraf, King-1437 shuk Dasgupta, Rasmus Larsen, Yicheng Wang, Man-1438 ish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki 1439 Osawa, Celine Smith, Ramya Sree Boppana, Tay-1440 lan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun, 1441 Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam 1442 Choo, Olaf Ronneberger, Chimezie Iwuanyanwu, 1443 Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene 1444 Cai, Shariq Iqbal, Martin Sundermeyer, Zhe Chen, 1445 Elie Bursztein, Chaitanya Malaviya, Fadi Biadsy, 1446 Prakash Shroff, Inderjit Dhillon, Tejasi Latkar, Chris 1447 Dyer, Hannah Forbes, Massimo Nicosia, Vitaly Niko-1448 laev, Somer Greene, Marin Georgiev, Pidong Wang, 1449 Nina Martin, Hanie Sedghi, John Zhang, Praseem 1450 Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Ji-1451 ageng Zhang, Viorica Patraucean, Dayou Du, Igor 1452 Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi 1453 Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan 1454 Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hud-1455 son, Vaishakh Keshava, Shubham Agrawal, Kevin 1456 Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Mad-1457 havi Sewak, Bryce Petrini, DongHyun Choi, Ivan 1458 Philips, Ziyue Wang, Ioana Bica, Ankush Garg, 1459 Jarek Wilkiewicz, Priyanka Agrawal, Xiaowei Li, 1460 Danhao Guo, Emily Xue, Naseer Shaik, Andrew 1461 Leach, Sadh MNM Khan, Julia Wiesinger, Sammy 1462 Jerome, Abhishek Chakladar, Alek Wenjiao Wang, 1463 Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Mar-1464 cus Wainwright, Mario Cortes, Frederick Liu, Joshua 1465 Maynez, Andreas Terzis, Pouya Samangouei, Ri-1466 ham Mansour, Tomasz Kepa, François-Xavier Aubet, 1467 Anton Algymr, Dan Banica, Agoston Weisz, An-1468 dras Orban, Alexandre Senges, Ewa Andrejczuk, 1469 Mark Geller, Niccolo Dal Santo, Valentin Anklin, 1470 Majd Al Merey, Martin Baeuml, Trevor Strohman, 1471 Junwen Bai, Slav Petrov, Yonghui Wu, Demis Has-1472 sabis, Koray Kavukcuoglu, Jeff Dean, and Oriol 1473 Vinyals. 2024. Gemini 1.5: Unlocking multimodal 1474 understanding across millions of tokens of context. 1475 Preprint, arXiv:2403.05530. 1476

Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun1477Woo, Manoj Middepogu, Sai Charitha Akula, Jihan1478Yang, Shusheng Yang, Adithya Iyer, Xichen Pan,1479Ziteng Wang, Rob Fergus, Yann LeCun, and Sain-
ing Xie. 2024a. Cambrian-1: A fully open, vision-1480

- IEEE/CVF Conference on Computer Vision and Pat-1482 centric exploration of multimodal llms. Preprint, 1539 arXiv:2406.16860. tern Recognition (CVPR), pages 9556–9567. 1483 1540 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, 1484 Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, and 1541 Yann LeCun, and Saining Xie. 2024b. Eyes wide Chen Change Loy. 2023. Contextual object detec-1485 1542 shut? exploring the visual shortcomings of multition with multimodal large language models. ArXiv, 1486 1543 1487 modal llms. Preprint, arXiv:2401.06209. abs/2305.18279. 1544 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, 1488 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, 1545 Yann LeCun, and Saining Xie. 2024c. Eyes wide 1489 and Lucas Beyer. 2023. Sigmoid loss for language 1546 shut? exploring the visual shortcomings of multiimage pre-training. *Preprint*, arXiv:2303.15343. 1547 modal llms. 2024 IEEE/CVF Conference on Com-1491 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun 1548 1492 puter Vision and Pattern Recognition (CVPR), pages Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, 1549 1493 9568-9578. Kai-Wei Chang, Peng Gao, and Hongsheng Li. 2024a. 1550 Guangzhi Wang, Yixiao Ge, Xiaohan Ding, Mohan S. 1494 Mathverse: Does your multi-modal llm truly see the 1551 1495 Kankanhalli, and Ying Shan. 2023. What makes for diagrams in visual math problems? In European 1552 1496 good visual tokenizers for large language models? Conference on Computer Vision. 1553 1497 ArXiv, abs/2305.12223. Zicheng Zhang, Haoning Wu, Erli Zhang, Guangtao 1554 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-1498 Zhai, and Weisi Lin. 2024b. Q-bench+: A bench-1555 hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin 1499 mark for multi-modal foundation models on low-1556 1500 Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei level vision from single images to pairs. Preprint, 1557 1501 Du, Xuancheng Ren, Rui Men, Dayiheng Liu, arXiv:2402.07116. 1558 1502 Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model's per-1503 ception of the world at any resolution. Preprint, 1504 arXiv:2409.12191. 1505 1506 Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu 1508
 - Sun, Qiong Yan, Guangtao Zhai, and Weisi Lin. 2024a. Q-bench: A benchmark for general-purpose foundation models on low-level vision. In The Twelfth International Conference on Learning Representations.
 - Penghao Wu and Saining Xie. 2023. V*: Guided visual search as a core mechanism in multimodal llms. 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 13084–13094.

1513

1515

1516

1517 1518

1519

1520

1521 1522

1523

1524

1526

1528

1529

1530 1531

1532

1533

1534

1535

1536

1537

- Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. 2024b. Deepseekvl2: Mixture-of-experts vision-language models for advanced multimodal understanding. Preprint, arXiv:2412.10302.
- Yijia Xiao, Edward Sun, Tianyu Liu, and Wei Wang. 2024. Logicvista: Multimodal llm logical reasoning benchmark in visual contexts. ArXiv. abs/2407.04973.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. 2024

A Detailed Related Work

Multimodal Large Language Models (MLLMs) have evolved from joint image-language pretraining approaches, exemplified by CLIP (Radford et al., 2021), BLIP (Li et al., 2022, 2023), and SigLIP (Zhai et al., 2023). With the rapid progress of large language models (LLMs), the 'vision as a token' paradigm has emerged, fostering numerous closed and open-source MLLMs that continue to advance through large-scale data and instruction finetuning (OpenAI et al., 2024b; Wu et al., 2024b; Wang et al., 2024). Various domain-specific benchmarks have been introduced in mathematics (Lu et al., 2023; Zhang et al., 2024a), science (Lu et al., 2022; Saikh et al., 2022), scientific documents (Li et al., 2024), and broader AGI contexts (Yue et al., 2023), each requiring intricate reasoning and domain knowledge. However, studies (Zhang et al., 2024a; Wu et al., 2024a) indicate that many VLM failures arise from inadequate visual perception, underscoring the need for robust, perception-focused evaluation.

Several datasets (Peng et al., 2023; Wang et al., 2023; Wu and Xie, 2023; Zang et al., 2023; Tong et al., 2024c; Wu et al., 2024a; Zhang et al., 2024b) do assess visual perception but often incorporate domain-specific reasoning. In contrast, visual benchmarks like MMVP (Tong et al., 2024b) and CV-Bench (Tong et al., 2024a) do provide insights into general visual perception abilities in MLLMs. However, these studies heavily rely on human annotation for dataset construction, which makes it difficult to scale them to larger sizes. Additionally, these datasets lack a mechanism to attribute difficulty levels to the visual perception of their constituent samples, which could be useful in determining fine-grained visual perception differences between MLLMs.

Insights from human psychology reveal that, visual perception encompasses five core dimensions (Chalfant and Scheffelin, 1969). Motor-free tests (Colarusso, 2003; Gardner, 1988), such as the Test of Visual Perceptual Skills (TVPS) (Gardner, 1988), Motor-Free Visual Perception Test (MVPT) (Colarusso, 2003), and Developmental Test of Visual Perception (DVPT) (Hammill et al., 2016) focus exclusively on perception, avoiding motor components. These, tests have evolved to assess children and adults (Brown and Peres, 2018), and have also been adapted for evaluating individuals recovering from visual cortex injuries (Brown and Peres, 2018).

Drawing on the extensive research in human psychology on human visual perception, we introduce the **Do You See Me** benchmark, leveraging established perception categories to evaluate VLMs.

B Do You See Me - Additional Details

B.1 Visual Discrimination

1587

1589

1590

1591

1592

1593

1594

1596

1598

1600

1601

1602

1604

1605

1607

1608

1610 1611

1612

1613

1614

1615

1616

1617

1619 1620

1621

1622

1625

The visual discrimination dimension of our dataset has four subtasks: shape discrimination, joint shapecolor discrimination, letter discrimination, and form-constancy. Below, we briefly outline our automated test generation for all the four visual discrimination sub-division.

Shape Discrimination: The difficulty of this task is controlled using: (a) number of unique shapes in the canvas (S), (b) maximum permissible instances per-shape (S_I) , and (c) allowable overlap factor (α) . The overlap factor α , refers to the separation amount for between shapes. When set to a non-negative value it results into a cluster of non-overlapping shapes, whereas, a negative value allows for some degree of overlap.

Letter Discrimination: The visual perception task's difficulty is modulated by three key factors: block spacing, letter count, and color contrast. The block spacing factor β determines the distance between a letter's constituent blocks, with $\beta = 0.1$ representing large spacing that makes letter identification challenging. As β decreases below 0.1, letter recognition becomes progressively easier. The task difficulty also increases with the number of letters present in the canvas. Color contrast provides another dimension of complexity, defined as $\Delta C = F - B$, where F represents the block color and B represents the background color. We implement three levels of contrast ($\Delta C \in 1, 2, 3$), where 1 indicates high foreground-background contrast and 3 represents low contrast. To generate these color combinations, we begin with seven common background colors in hexadecimal format, converting them to HSV representation. We then create three difficulty levels by manipulating the hue, saturation, and value: for easy combinations (high contrast), we rotate the hue by 180, for medium combinations, we rotate the hue by 90 while maintaining constant saturation and value, for challenging combinations (low contrast), we rotate the hue by 30 and apply a 10% distortion to both saturation and value. These parameters are summarized in Table 4

Visual Form Constancy: It is important to note that low rotation and scaling without substitution make the task more challenging, whereas higher values of rotation, scaling, or substitution simplify it by rendering the variants more visibly distinct.

B.2 Visual Closure

We begin by defining seven basic shapes: capsule, star, hexagon, circle, pentagon, rectangle, and triangle. Each shape is represented by a set of vertices connected to form the exact shape. A single shape is randomly selected as the target image. To create an incomplete version of this shape, we first choose k pairs of adjacent vertices and remove the connecting edges entirely. Next, we choose l other pairs of adjacent vertices (distinct from the first k pairs) for partial edge removal. This produces an incomplete shape which, if completed, would match the original target. To generate three distractors (incorrect options), we take the incomplete shape, introduce distortions to m vertices by adding noise specified by the distortion factor δ to the vertex coordinates. Each of the three distractors is created by repeating this process independently, resulting in three unique noisy variants. Finally, the correct incomplete shape (which can be closed to match the target) and the three distractors are shuffled and displayed horizontally, accompanied by the complete target shape on the top (see Fig 3).

C Detailed Results

Table 2: Performance comparison across different visual perception tasks. All values are percentages (%). Models are sorted by average accuracy.

Model	Figure	Visual	Color	Shape	Letter	Form	Visual	Average
	Ground	Spatial	Disamb.	Disamb.	Disamb.	Const.	Closure	Accuracy
Human	100.00	92.59	100.00	100.00	77.77	98.14	91.66	94.31
Claude Sonnet-3.5	47.78	32.63	75.98	30.83	13.33	91.48	58.33	50.05
Gemini-1.5	34.44	25.68	81.86	26.67	22.96	81.11	67.26	48.57
GPT-40	33.33	25.43	66.91	10.83	28.89	80.37	58.93	43.53
Qwen-2.5	27.78	40.69	81.86	19.58	2.96	50.37	63.10	40.91
Deepseek-tiny	20.00	21.34	57.35	20.42	25.19	30.74	29.76	29.26
Intern-VL	30.00	16.00	55.64	19.17	0.74	34.81	29.76	26.59
Deepseek-small	27.78	14.52	45.10	10.00	1.48	43.70	35.71	25.47
LLaMA-11B	26.67	6.08	25.74	3.33	11.85	22.22	22.02	16.84
LLaMA-90B	23.33	10.05	15.44	1.25	11.11	27.41	21.43	15.72

Table 3: Comparison of model performance. Claude Sonnet-3.5 leads in both *reasoning* and *visual perception* questions.

Model	Reasoning Acc. (%)	Perception Acc. (%)
Claude Sonnet-3.5	40.95	45.21
GPT-40	32.97	42.55
Gemini 1.5 Flash	32.97	44.68
Qwen2.5-VL-7B-Instruct	35.63	35.10
Intern2.5-VL-8B	27.66	37.23
LLaMA3.2-Vision-Preview-90B	35.64	25.00
LLaMA3.2-Vision-Preview-11B	26.06	31.91
Deepseek-VL-Small 3B	22.34	26.60
Deepseek-VL-Tiny 1B	23.40	35.01



Figure 9: MLLM performance on Human Rated Difficulty Levels. Note: *Empty human bar for a difficulty level indicates that no samples were attributed the corresponding difficulty level.*



Figure 10: Average MLLM performance over a sweep of combinations of control parameters.

D Do You See Me - Control Parameter Details

1627

Division	Subdivision	Control Parameters	Question Type	Unique Images	Number of Questions
	Shape Discrimination	Number of Shapes: $S \in [3, 7]$ Instances per Shape: $S_I \in [3, 6, 10]$ Overlap Factor: $\alpha \in [-40, -30, -20, 10]$	Integer	241	241
	Joint Shape-Color	Number of Shapes: $S \in [2, 4, 6]$ Number of Unique Colors: $C \in [2, 4, 6]$	Integer	90	408
Visual Discrimination	Letter Discrimination	Number of Letters: $N \in [1, 5, 9]$ Foreground-Background Contrast: $\Delta C \in [1, 2, 3]$ Block Size: [0.04, 0.08, 0.1]	Text	135	135
	Form Constancy	Shape Substitution Factor : $ssf \in [0, 1]$ Scaling Factor: $\alpha \in [0.8, 1.1, 1.4]$ Rotation Factor: $\theta_r \in [5, 25, 50]$ Aspect Ratio: $\beta \in [0.8, 1.1, 1.4]$	MCQ	270	270
Visual Spatial	Spatial Grids	Grid Dimension: $D \subseteq [3, 6, 9] \times [3, 6, 9]$ Number of Grids: $G \in [1, 3, 5]$	Integer	270	806
Visual Figure-Ground	N.A.	Number of Shapes: $N \in [2, 6, 10]$ Background Density Factor: $bdf \in [0.1, 0.3, 0.5]$	MCQ	90	90
Visual Closure	N.A.	Number of Full Edges to Remove: $k \in [1, 3]$ Number of Partial Edges to Remove: $l \in [1, 3]$ Number of Edges to Distort: $m \in [1, 3]$ Distortion Factor: $\delta \in [0.1, 0.12, 0.14]$	MCQ	166	166

Table 4: Control parameters and question types for each subdivision of the visual perception test.

E Human Performance Benchmarking

We recruited 7 participants, consisting of 5 men and 2 women. The selection criterion was that participants1629had no history of vision-related disorders. All participants were between 25 and 35 years of age. Each test1630took approximately 10 minutes to complete, with a total testing time of 70 minutes. To prevent mental1631fatigue, we conducted the tests over two days.1632

Welcome to Visual Perception Test

Enter your name	
Select your gender	
Male	~
Enter your age	
1 -	+
Select test type	
geometric_dataset	~
Start Test	

(a) Information Collected in the Test

Test Instructions -Geometric_dataset Test

- 1. You will first go through a calibration phase with 7 practice questions
- 2. After calibration, you will proceed to the actual test
- 3. Each question has a timer
- 4. Enter your answer in the text box provided
- 5. Click 'Submit' to move to the next question
- 6. After each question in the actual test, rate its difficulty level

Start Calibration

(b) Calibration Instructions





Figure 12: Actual Test and Human Difficulty Rating Collection

F Joint Perception-Reasoning Dataset

F.1 Correct Reasoning Incorrect Perception

Example 1



Figure 13

Reasoning Question: Which of the boxes comes next in the sequence? Select answers from A-E. Explain the reasoning behind the answer.

Reasoning Answer (GT): C

Reasoning Explanation: Triangle alternates from pointing up, to pointing right with each turn. Doughnut shape alternates from left corner to right corner with each turn.

Claude Sonnet-3.5 Answer: Let's solve this step by step:

Looking at the top sequence, I notice a pattern of shapes and dots:

Single triangle, single dot above Single triangle, single dot above Single triangle, two dots above Single triangle, single dot above Single triangle, single dot above Also notice that the triangles alternate between pointing left () and right ()

The pattern appears to be:

The triangles alternate direction The dots follow a pattern of: 1, 1, 2, 1, 1, 2, ... Following this pattern, the next box should have:

A triangle pointing right () Two dots above (because we're at the point in the sequence where two dots should appear) Looking at the options A through E: Option C shows a triangle pointing right with two dots above it.

Therefore, the answer is C.

The reasoning is based on both the alternating direction of the triangles and the repeating pattern of dots (1, 1, 2).

VLM Final Answer: YES

Visual Perception Question: Using visual inspection, count the number of black dots on the top-left side of the box where the triangle present in the box is resting on its side and not the vertex. **Visual Perception Answer (GT):** 4

Claude Sonnet-3.5 Answer: Let me help you count the black dots as requested.

Looking at the boxes where the triangle is resting on its side (not on its vertex), we need to look at the top-left portion of those boxes.

The triangles resting on their sides (not on vertex) appear in positions B, D, and E in the top row and positions A and E in the bottom row.

Among these:

Top row: B has 2 dots, D has 2 dots, E has 1 dot Bottom row: A has 0 dots, E has 2 dots (**incorrect perception**) Therefore, there are a total of 7 black dots when counting only the top-left portions of boxes where the triangle is resting on its side.

Analysis: The visual perception question shows how model fails answering the question on number of dots on top-left side. The same failure occurs in answering the visual reasoning question.

- -



Figure 14

Reasoning Question: Which set does the Figure belong to? Select from A, B, and C. (A) Set A (B) Set B (C) Neither set A nor set B. Explain the reasoning behind the answer.

Reasoning Answer (GT): B

Reasoning Explanation: Set A: The total number of crossed boxes is equal to the total number of stars. The total number of shaded boxes is equal to the total number of circles. Set B: The total number of crossed boxes is equal to the total number of circles. The total number of shaded boxes is equal to the total number of starts. Figure: There are four shaded boxes and four stars. There are three crossed boxes and three circles.

Claude Sonnet-3.5 Answer: The Figure belongs to (B) Set B.

Reasoning:

Looking at the pattern in both sets, we need to focus on the arrangement of black squares and circles in the cards.

Key observations about the Figure:

It has circles arranged in a specific pattern It has three black squares at the bottom The overall layout is similar to other cards in both sets When comparing with Set A and Set B: Set A's cards generally have black squares distributed throughout different positions Set B's cards consistently show black squares grouped at the bottom of each card, similar to the Figure The circle patterns in Set B also appear more similar to the Figure's arrangement The decisive factor is the positioning of the black squares: The Figure shows three black squares aligned at the bottom This arrangement is consistent with the pattern seen in Set B, where black squares are consistently placed at the bottom of the cards This pattern is not characteristic of Set A Therefore, based on the consistent pattern of black square placement at the bottom and similar circle arrangements, the Figure belongs to Set B. **MLLM Final Answer:** Correct

Visual Perception Question: In Set B of the diagram, looking at the second square from the left, how many stars (?) are inside that square?

Visual Perception Answer (GT): 4

Claude Sonnet-3.5 Answer: Looking at Set B, in the second square from the left, I can count 2 stars within that square.

MLLM Final Answer: Incorrect

Analysis: The that number of stars within the box are linked to correctly the reasoning question and Claude Sonnet-3.5 fails to answer the count correctly. Further, from the reasoning chain, it is clear that it makes a reasoning error that black squares are grouped/ungrouped in Set A and B, and answers the reasoning question on this basis. This shows a reasoning error in its response, however, the final answer given is correct. Additionally, there is visual hallucination about distribution of black squares as well.

F.2 Qualitative Analysis of Reasoning Chain Errors





Figure 15

Reasoning Question: Which of the boxes comes next in the sequence? Select from A-E **Reasoning Answer (GT):** A

Reasoning Explanation:

GPT4o Answer: To determine which box comes next in the sequence, observe the direction of the arrows and the distribution of dots in each box. 1. The direction of arrows alternates from pointing up to down in every subsequent box. The sequence starts with an arrow pointing up. 2. The number of dots increases incrementally by one in each subsequent box. The first box has one dot, the second has two, the third has three, and so on. Following this pattern:

- The first box in the sequence has one dot and an arrow pointing up.

- The second box has two dots and an arrow pointing down.

- The third box has three dots and an arrow pointing right.

- The fourth box has four dots and an arrow pointing down. (perception error)

- The fifth box has five dots and an arrow pointing up.

Given these patterns, the next box should have an arrow pointing down and six dots. Among the choices A to E:

- Choice A fits with an arrow pointing down and six dots.

Therefore, the correct answer is box A.

VLM Final Answer: Correct.

Key Observation: Correct Final Answer But presence of Incorrect Visual Perception in Reasoning Chain.

Analysis: GPT40 makes a visual perception error on the fourth box, the pattern was up, down, right, left, up. Whereas, GPT-40 determines that the pattern is up, down, right, down, up. *This is an example of visual perception error in reasoning chain even while the model gives a correct final answer.*



Figure 16

Reasoning Question: Which of the boxes comes next in the sequence? Select from A-E **Reasoning Answer (GT):** C

Reasoning Explanation:

GPT40 Answer: To determine which box comes next in the sequence, we need to observe the patterns of the shapes and their rotations in the given boxes. 1. The position of the circle:

- The circle moves through four quadrants: top-right, top-left, bottom-left, and bottom-right. (visual perception error: the circle never goes to bottom left or right)

- Based on this pattern, after the **bottom-right** position in the fifth box, the next position should be top-right.

2. The rotation of the triangle:

- The triangle rotates clockwise: pointing up, right, down, and then left.

- In the fifth box, the triangle is pointing left. Therefore, in the next box, it should point up again.

Given these patterns, we are looking for a box with the circle in the top-right quadrant and the triangle pointing up.

Option $**C^{**}$ shows a circle in the top-right quadrant and a triangle pointing up. This matches the patterns observed, so option C is the correct answer for the next box in the sequence.

VLM Final Answer: Correct.

Key Observation: Correct final answer but *presence of incorrect visual perception* in Reasoning Chain.

Analysis: The circle goes between top-right/left. Example of visual perception error

•					?
$\stackrel{\wedge}{\rightarrow} \stackrel{\wedge}{\rightarrow}$	*	* * *	*	$\stackrel{\diamond}{\bullet}$	$\uparrow \uparrow \uparrow$

Figure 17

Reasoning Question: The bottom boxes create a rule that has to be applied in the box directly above them. Select which of options A to F corresponds to the rule below the box with the question mark.. Explain the reasoning behind the answer.

Reasoning Answer (GT): E

Reasoning Explanation: The simple rule in this question is that the number of grey arrows in the box below must equal the number of black shaded shapes in the box above. There are three grey arrows below the missing box, so the correct answer must contain three black shapes.

GPT40 Answer: To determine which option (A to F) corresponds to the rule that should be applied above the question mark, we need to analyze the pattern in each column and how the bottom boxes affect the top boxes.

Column Analysis:

Column 1: The bottom box has arrows pointing right. The top box changes the orientation of the shapes accordingly.

Column 2: The bottom box has arrows pointing left. The top box changes by mirroring or flipping the orientation.

Column 3: The bottom box has inward-facing arrows, and the shapes in the top box are transformed or interchanged.

Column 4: The bottom box has outward-facing arrows. It seems to involve swapping the positions of certain elements in the top box.

Column 5 (Question Mark): The bottom box has left-facing arrows, implying a mirroring or flipping operation similar to column 2.

Rule Application:

Given the bottom box under the question mark has left-facing arrows, we should select an option that mirrors or flips the elements horizontally compared to the option from Column 4.

Options Evaluation:

Option E is a horizontally flipped version of the top box in Column 4, making it the correct answer. Therefore, the correct option corresponding to the rule below the box with the question mark is Option E.

VLM Final Answer: Correct

Key Observation: Correct answer, incorrect reasoning.

Analysis: The simple rule here is the number of grey arrows in the box below must equal the number of black shaded shapes in the box above. *However, the model reasons about changing orientations based on the direction of arrows.* However, even with wrong reasoning GPT-40 reaches the correct answer.



Figure 18

Reasoning Question: Which set does the Figure belong to? Select from A, B, and C. (A) Set A (B) Set B (C) Neither set A nor set B. Explain the reasoning behind the answer.

Reasoning Answer (GT): C

Reasoning Explanation: Set A: Each box contains four shaded shapes. Set B: Each box contains three shaded shapes. Figure: The Figure contains five shaded shapes.

GPT4o Answer: The figure belongs to (C) Neither set A nor set B. Reasoning: Set A:

In Set A, the figures alternate between filled and empty shapes, or there is a symmetry of arrangement. Set B:

In Set B, there is a symmetrical arrangement as well, but the colors (filled and empty) alternate in a different manner than Set A. Figure Analysis:

The given figure is composed of triangles, both filled and empty, arranged in a different pattern from any of the figures in Set A or Set B. It does not follow the alternation pattern or symmetry rules observed in either set. Therefore, the figure does not belong to Set A or Set B.

VLM Final Answer: Correct

Key Observation: Correct answer, incorrect reasoning.

Analysis: GPT4o reasons that there is a different alternate pattern to fill shapes in Set A and B, and the new shape has both these patterns, hence cannot be alloted to any one set. However, the simple rule here is that 3 shapes are shaded in set B, and 4 shapes are shaded in shape A. Thus, while the model reasons incorrectly, it arrives at the right option.

G MLLM Prompts

Expert LLM Evaluator

You will be given three inputs consisting of: a question, the answer by a human giving the test, and the ground truth answer. The answer from the human can be a longform one, which can contain an explaination or reasoning. Your task is to determine if the final answer reached matches the ground truth answer.

Inputs

Question: {question}

Ground Truth Answer: {gt_answer}

Human Answer: {mllm_answer}

Your task is to reply in just YES (or) NO. If the answer matches your response should be YES, else NO.

1643

1642

Example Prompts For Subtasks in Do You See Me

Visual Figure Ground: The figure consists of a Target image, which is embedded in some background noise. Out of the four given options, your task is to pick the option which has the same figure as the target image. Respond as follows: Option your answer (choose between 1, 2, 3, or 4).

Letter Disambiguation: The image shows one or more letters formed by a grid of small squares. What letter(s) can you identify in this image? Please respond with only the letter(s) you see.

Visual Form Constancy: The figure consists of a Target image. Out of the four given options, your task is to pick the option which has the same figure as the target image. Respond as follows: Option your answer (choose between 1, 2, 3, or 4).

Visual Closure: The figure consists of a target image which is complete, Out of the four given options (which are partially complete), your task is to pick the option which when completed matches the target image. Respond as follows: Option your answer (choose between 1, 2, 3, or 4).

Visual Spatial: In grid 5, starting from the white square at position (row 1, column 5), how many circles are there down of it in the same column?

Color Disambiguation: Count the number of cross's that are purple.

Shape Discrimination: Count the total number of stars in the image, including each concentric star separately. For example, if there is one star with 2 inner concentric rings, that counts as 3 stars. Respond with only a number.