Jigsaw-Puzzles: From Seeing to Understanding to Reasoning in Vision-Language Models

Anonymous ACL submission

1

Abstract

Spatial reasoning is a core component of human cognition, enabling individuals to perceive, comprehend, and interact with the physical world. It relies on a nuanced understanding of spatial structures and inter-object relationships, serving as the foundation for complex reasoning and decision-making. To investigate whether current vision-language models (VLMs) exhibit similar capability, we introduce Jigsaw-Puzzles, a novel benchmark consisting of 1,100 carefully curated real-world images with high spatial complexity. Based on this dataset, we design five tasks to rigorously evaluate VLMs' spatial perception, structural understanding, and reasoning capabilities, while deliberately minimizing reliance on domainspecific knowledge to better isolate and assess the general spatial reasoning capability. We conduct a comprehensive evaluation across 24 state-of-the-art VLMs. The results show that even the strongest model, Gemini-2.5-Pro, achieves only 77.14% overall accuracy and performs particularly poorly on the Order Generation task, with only 30.00% accuracy, far below the 90%+ performance achieved by human participants. This persistent gap underscores the need for continued progress, positioning Jigsaw-Puzzles as a challenging and diagnostic benchmark for advancing spatial reasoning research in VLMs.

1 Introduction

011

017

042

The road to artificial general intelligence (AGI) demands more than language or vision alone: it requires models to possess a human-like spatial reasoning capability by constructing structured representations of the physical world (Lake et al., 2017). Spatial reasoning refers not just to the perception of visual input, but to the capability to comprehend spatial arrangements, model structural relations, and infer the geometry and layout of a scene. These capabilities are fundamental to human cognition



Figure 1: Jigsaw-Puzzles example. While human participants effortlessly reconstruct the original spatial layout, all tested VLMs fail to recover the correct order.



Figure 2: Evaluation of VLMs on Jigsaw-Puzzles. The plot reports the accuracy of 8 representative VLMs on 5 tasks.

and develop naturally through everyday perception and interaction (Ishikawa and Newcombe, 2021). In contrast, current VLMs, while highly capable in tasks such as image captioning (Young et al., 2014; Lin et al., 2014; Sharma et al., 2018), visual question answering (Krishna et al., 2017; Singh et al., 2019; Marino et al., 2019), and image-text retrieval (Schuhmann et al., 2021; Thapliyal et al., 2022; Bitton-Guetta et al., 2023), consistently struggle with tasks requiring spatial reasoning (Stogiannidis

043

100

101

102

104

et al., 2025). We show an example in Figure 1 and report the performance of some tested VLMs on Jigsaw-Puzzles in Figure 2. This gap underscores a critical limitation: while current VLMs have made substantial progress in basic visual understanding, they continue to struggle with structured spatial reasoning, which is essential for grounded understanding in real-world scenarios. Bridging this gap is essential for progressing towards generalizable human-like spatial cognition and ultimately AGI.

However, existing benchmarks have yet to provide a comprehensive evaluation of spatial reasoning capability in VLMs under complex, realworld visual scenarios. Some works (Fu et al., 2024; Li et al., 2024; Liu et al., 2024; Yue et al., 2024) focus primarily on foundational visual understanding by systematically evaluating perception, comprehension, and basic visual reasoning, revealing notable limitations in these areas. Although a few recent efforts (Pothiraj et al., 2025; Stogiannidis et al., 2025; Ren et al., 2025; Tang et al., 2025) have attempted to evaluate the spatial reasoning capability of VLMs, they often rely on overly synthetic settings, task-specific constraints, or domain-dependent priors (Song et al., 2025) (See Appendix A.1 for examples), limiting the capability to capture generalizable spatial reasoning under natural visual conditions. A truly effective evaluation of human-like spatial reasoning capability should model the task as a multi-stage cognitive process-beginning with perception, advancing through structural understanding, and culminating in high-level reasoning. Such reasoning must be grounded in the visual richness and ambiguity of real images, requiring the integration of spatial structural modeling and goal-directed reasoning (Chen et al., 2024). Yet, this critical dimension of spatial cognition remains largely overlooked in existing benchmarks, underscoring the need for new benchmarks that move beyond narrow task formulations and embrace the full complexity of spatial reasoning.

To overcome the limitations of existing benchmarks in evaluating spatial reasoning capability of VLMs under real-world conditions, we introduce Jigsaw-Puzzles, a novel benchmark inspired by the cognitive mechanisms underlying human puzzlesolving. Puzzle-solving naturally reflects the multistage cognitive process (Fissler et al., 2018), making it a compelling and effective testbed for spatial reasoning in VLMs.

In total, Jigsaw-Puzzles comprises 1,100 care-

Bonohmowk	Understanding	Beaconing	High Visual	Great	Automated	
Benchinai K	Understanding	Reasoning	Complexity	Scalability	Construction	
Capture (Pothiraj et al., 2025)	~	~	×	×	×	
Mind the Gap (Stogiannidis et al., 2025)	~	~	×	×	×	
VGRP (Ren et al., 2025)	√	~	×	×	×	
LEGO-Puzzles (Tang et al., 2025)	~	~	×	~	√	
Jigsaw-Puzzles (Ours)	~	~	~	~	1	

Table 1: Comparison of spatial reasoning benchmarks.

105

106

107

108

109

110

111

112

113

114

115

116

117

118

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

150

fully curated real-world images and features five different tasks. First, we begin with the *Missing Piece Selection* task to evaluate VLMs' basic spatial understanding capability, which serves as the essential foundation for spatial reasoning. Building on this foundation, we introduce four reasoningcentric tasks: *Piece Localization, Connection Verification, Anomaly Detection*, and *Order Restoration.* These tasks are designed to assess various facets of spatial reasoning, including adjacency modeling, local structural consistency, spatial localization, geometric transformation understanding, and multi-step spatial reasoning.

Compared to existing benchmarks for evaluating spatial reasoning capability in VLMs, Jigsaw-Puzzles offers three key advantages, as summarized in Table 1: (1) Higher visual complexity. Jigsaw-Puzzles uses real-world images with diverse and rich visual elements, and significantly outperforms benchmarks based on synthetic images (Ren et al., 2025; Stogiannidis et al., 2025; Tang et al., 2025) and simple visual scenes (Pothiraj et al., 2025). This enables more realistic and challenging spatial reasoning evaluation. (2) Greater scalability. Any natural image that satisfies the construction rules can be directly used to generate puzzle tasks, without the need to manually synthesize target images. (3) Fully automated construction pipeline. All Jigsaw-Puzzles tasks are generated automatically without manual annotation, with each question paired with a unique deterministic answer. This feature enables low-cost dataset construction and facilitates continuous expansion and refinement.

In summary, we introduce Jigsaw-Puzzles, a novel benchmark for systematically evaluating the human-like spatial reasoning capability of VLMs in realistic visual settings. Our main contributions are as follows:

A new benchmark for spatial reasoning. We introduce Jigsaw-Puzzles, a puzzle-inspired benchmark constructed through a fully automated pipeline that improves existing benchmarks in visual complexity and scalability, while enabling structured evaluation of spatial reasoning in VLMs. **Comprehensive evaluation and analysis.** We evaluate 24 state-of-the-art VLMs on Jigsaw-

Puzzles and conduct detailed analysis. Our findings 151 expose consistent limitations in current VLMs and 152 provide insights to guide future improvements in 153 spatial reasoning capability. 154

Open-sourced dataset and construction tools. 155 We release the full dataset along with the auto-156 mated generation scripts to support the evaluation 157 and continued advancement of spatial reasoning in 158 VLMs under real-world visual scenarios, as well 159 160 as to facilitate future benchmark expansion.

Related Work 2

161

163

164

165

167

168

169

170

171

173

174

175

177

179

180

181

182

184

185

187

General VLMs Evaluation Benchmarks. With the rapid progress of VLMs, systematically evaluating their diverse capabilities has become a key challenge. Although many benchmarks have been introduced, most focus primarily on visual understanding. MME (Fu et al., 2024) evaluates instruction following, perception, and basic reasoning across 14 subtasks, revealing persistent issues like object hallucination and limited spatial understanding. SEED-Bench (Li et al., 2024) includes 19,000 multiple-choice questions across 12 dimensions 172 and shows continued struggles with text recognition and temporal reasoning. MMBench (Liu et al., 2024) offers fine-grained bilingual evaluations, enhancing the robustness of multilingual assessment. 176 MMMU (Yue et al., 2024) provides 11,500 questions across 183 subfields and 30 image types to 178 test expert-level reasoning, yet even advanced models like Gemini display notable knowledge gaps. While these benchmarks have advanced the evaluation of perceptual and semantic understanding, none systematically assess spatial reasoning-the core aspect of human cognition. This highlights the pressing need for more challenging and diagnostic benchmarks specifically targeting spatial reasoning capability in VLMs.

Spatial Reasoning Evaluation Benchmarks in 188 VLMs. Several recent benchmarks have aimed to 189 evaluate the spatial reasoning capability of VLMs. 190 Capture (Pothiraj et al., 2025) assesses occluded object counting, revealing that VLMs struggle to 192 form coherent spatial representations under occlu-193 sion. Mind the Gap (Stogiannidis et al., 2025) 194 evaluates spatial relations, navigation, and men-195 196 tal rotation, showing that VLMs often perform near chance level, indicating limited spatial cognition. 197 VGRP (Ren et al., 2025) introduces 20 visual grid 198 puzzles across varying difficulty levels to assess visual perception, rule-following, and logical rea-200

soning. LEGO-Puzzles (Tang et al., 2025) provides 1,100 visual QA pairs over 11 subtasks to measure basic and multi-step spatial reasoning. Results consistently show that current VLMs struggle with perceptual complexity, rotation reasoning, and sequential reasoning. Despite these efforts, most benchmarks rely on simplified scenarios, failing to reflect the complexity of real-world spatial environments, thereby limiting their generalizability. More diagnostic benchmarks grounded in natural visual settings are needed to advance human-level spatial reasoning in VLMs.

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

224

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

Jigsaw-Puzzles 3

In this section, we introduce Jigsaw-Puzzles, a scalable and comprehensive benchmark designed to evaluate the spatial reasoning capability of VLMs in realistic visual environments. Specifically, Section 3.1 outlines the motivation and definition of each task, while Section 3.2 describes the dataset construction process, including image selection and the automated generation of question-answer pairs.

3.1 **Task Definition**

To systematically evaluate the spatial reasoning capability of VLMs, we design tasks around the core cognitive stages underlying human spatial reasoning-seeing, understanding, and reasoning. Inspired by the human process of solving jigsaw puzzles, our benchmark simulates how individuals integrate fragmented visual information into a coherent whole: beginning with the perception of local visual cues, followed by the comprehension of spatial relationships, and culminating in multi-step spatial reasoning to reconstruct the original scene. This sequence naturally reflects the progression from low-level perception to high-level spatial reasoning. Accordingly, the tasks span spatial understanding, single-step, and multi-step spatial reasoning, collectively providing a comprehensive evaluation across different levels of spatial reasoning. Figure 3 shows examples of each task in Jigsaw-Puzzles.

Task 1: Missing Piece Selection. The task evaluates VLMs' spatial understanding capability. Given an image with a missing region, VLMs need to select the correct patch from four candidates. We define two difficulty levels: Easy, where distractors are randomly chosen, and Hard, where distractors are selected using CLIP (Radford et al., 2021) similarity to closely resemble the ground-truth patch, increasing the task's difficulty. Task 2: Piece Lo-



Figure 3: Task examples of Jigsaw-Puzzles. Note: the questions above are slightly simplified for clarity and brevity, and the **blue** option indicates the correct answer.



Figure 4: Dataset curation pipeline. Step 1 filters candidate images through expert-defined rules to build a spatial reasoning dataset. Step 2 uses automated templates to generate task-specific QA pairs from the curated images.

calization. This task evaluates spatial localization as a representative single-step spatial reasoning capability. Given a partially masked image and one masked patch, VLMs must identify the patch's original position. Difficulty is controlled by grid size and number of masked regions: Easy (2×2 with two masks), Hard (3×3 with four masks), increasing spatial complexity. *Task 3: Connection Verification.* This task evaluates adjacency reasoning, which also falls under single-step spatial rea-

256

soning. The full image is divided into 2×2 grids, and two patches are randomly selected. VLMs are asked to determine their spatial relationship in the original image (e.g., above-below, left-right, or non-adjacent). *Task 4: Anomaly Detection*. This task targets local spatial transformation detection, a process that inherently involves single-step spatial reasoning. One region in 2×2 grids is randomly rotated, mirrored, or left unchanged. The model must detect the change, locate the region, and iden-

tify the transformation. *Task 5: Order Restoration.*This task integrates the capabilities assessed in the
previous tasks and serves as a complex multi-step
spatial reasoning challenge. A complete image is
split into four shuffled patches. VLMs must identify the correct order to reconstruct the original
spatial layout.

Overall, the five puzzle-inspired tasks in Jigsaw-Puzzles cover a broad spectrum of spatial reasoning challenges—from basic spatial understanding to single-step and multi-step spatial reasoning—enabling a comprehensive evaluation of spatial reasoning in VLMs.

3.2 Dataset Curation

278

279

284

286

290

291

296

297

298

299

301

302

303

305

308

310

As illustrated in Figure 4, our dataset curation pipeline consists of two main stages: data collection and QA generation.

Data Collection. We integrate data collection and quality control into a unified process. Starting from the CC3M (Sharma et al., 2018) dataset, we apply task-specific filtering criteria—including minimum resolution and aspect ratio constraints—to construct an initial image pool of approximately 10,000 candidate images. Two human experts iteratively review the image pool while incrementally refining a shared set of filtering rules. Based on these evolving rules, they collaboratively filter the initial dataset to obtain the final set of high-quality, structurally diverse images. See Appendix A.2 for the rules pool. To enhance generalizability, we emphasize semantic and structural diversity throughout the dataset.

QA Generation. To support scalable and consistent QA pairs generation, each task type is associated with a specific construction template. QA pairs are automatically generated using these templates. Figure 3 illustrates simplified examples of the templates, full versions are provided in Appendix A.2.

4 Evaluation on Jigsaw-Puzzles

4.1 Experimental Setting

Benchmark Models. We evaluate 24 VLMs 311 on Jigsaw-Puzzles, covering a diverse range 312 of model scales and training paradigms. For 313 314 open-source models, we evaluate Qwen2-VL-72B (Wang et al., 2024a), QvQ-72B-Preview 315 (Qwen, 2024), Qwen2.5-VL-[7B/32B/72B] (Bai et al., 2025), InternVL3-[8B/14B/38B/78B] (Zhu et al., 2025), Kimi-VL-A3B-[Instruct/Thinking] 318

(Du et al., 2025), Phi-4-multimodal-instruct (Abouelenin et al., 2025), Aya-Vision-[8B/32B] (Dash et al., 2025), and Mistral-Small-3.1-24B-Instruct (Mistral, 2025). For proprietary models, we evaluate Claude-[3.5/3.7]-Sonnet (Anthropic, 2024), Gemini-[2.0/2.5]-Flash, Gemini-2.5-Flash-Thinking, Gemini-2.5-Pro (Anil et al., 2023), GPT-40, GPT-40-mini (Achiam et al., 2023), and Grok-2-Vision (Grok, 2024). Notably, QvQ-72B-Preview, Kimi-VL-A3B-Thinking, Gemini-2.5-Flash-Thinking, and Gemini-2.5-Pro are categorized as reasoning-enhanced models. All models, supporting multi-image input, are evaluated in a zero-shot setting with hardware scaled to their parameter size, see details in Appendix A.2.

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

354

357

358

359

360

361

362

363

364

365

366

367

368

369

Evaluation Metric. Since each QA pair in Jigsaw-Puzzles has a single correct answer, we use exact match accuracy (%) as the primary metric to evaluate VLMs' performance on each task.

Baselines. We provide two baselines for comparison: (1) Random, which assumes equal probability across all options and calculates expected accuracy accordingly. (2) p-value-based critical value, which reports the minimum accuracy required to outperform random guessing at a significance level of p=0.05.

Human Performance. To evaluate human performance, we construct a subset called Jigsaw-Puzzles-Lite by sampling 220 images from the full dataset. Three human participants complete all tasks on this subset under the same conditions as VLMs—without access to any external tools or the internet. Their performance serves as an empirical upper bound for spatial reasoning capability.

4.2 Main Results

Table 2, 3 report the performance of 24 VLMs on Jigsaw-Puzzles. Building on these results, we conduct a comprehensive and systematic analysis. We summarize several key findings as below.

Spatial Reasoning Remains a Challenge for VLMs. As shown in Table 3, human participants consistently outperform VLMs, achieving an overall accuracy of 96.36%. By comparison, current VLMs perform considerably worse, with even the strongest models—Gemini-2.5-Pro—lagging more than 20 percentage points behind human accuracy across all tasks. The persistent gap between humans and VLMs highlights the demanding nature of Jigsaw-Puzzles and affirms its utility as a robust benchmark for spatial reasoning evaluation. Significant Gap Between Open-Source and

Madala	Missing P	iece Selection	Piece Localization		Connection	onnection Anomaly		Overall
Widdels	Easy	Hard	Easy	Hard	Verification	Detection	Restoration	Overall
Baseline								
Random Guessing	25.00	25.00	50.00	25.00	33.33	28.13	25.00	30.21
\uparrow Random (p < 0.05)	27.30	27.30	52.50	27.30	35.70	30.50	27.30	32.56
Proprietary Models								
Grok2-Vision	64.55	52.45	53.00	41.00	34.91	27.73	25.27	42.70
GPT-4o-mini	96.45	83.64	59.45	37.82	44.36	57.91	33.18	58.97
GPT-4o	95.00	89.18	61.55	53.45	41.09	53.18	31.55	60.71
Claude-3.5-Sonnet	99.73	94.55	62.45	41.09	45.64	67.45	35.00	63.70
Claude-3.7-Sonnet	99.55	95.09	60.27	44.55	47.91	67.00	39.82	64.88
Gemini-2.0-Flash	92.09	85.64	63.55	54.00	44.91	68.73	34.27	63.31
Gemini-2.5-Flash	98.82	92.45	64.55	54.55	48.82	67.36	34.45	65.86
Gemini-2.5-Flash-Thinking	99.55	94.73	76.64	51.27	57.91	62.00	64.82	72.42
Gemini-2.5-Pro	99.91	97.18	78.82	61.09	59.36	70.00	73.64	77.14
Open-source Models								
Kimi-VL-A3B-Instruct	67.91	52.55	51.45	29.82	37.91	21.82	32.82	42.04
Kimi-VL-A3B-Thinking	84.64	58.09	56.36	32.64	28.00	30.91	20.36	44.43
Phi-4-multimodal-instruct	63.45	51.64	60.91	37.45	36.64	43.36	27.64	45.87
Qwen2.5-VL-7B	87.18	63.27	54.27	36.18	38.55	45.00	28.09	50.36
Aya-Vision-8B	26.82	27.27	49.64	26.55	35.00	12.91	24.73	28.99
InternVL3-8B	98.09	85.45	53.91	35.45	44.91	46.82	34.36	57.00
InternVL3-14B	99.73	88.73	59.64	40.18	51.73	49.09	40.91	61.43
Mistral-Small-3.1-24B-Instruct	26.91	28.55	54.27	31.27	38.27	52.36	26.45	36.87
Qwen2.5-VL-32B	97.27	76.82	62.09	40.36	50.09	61.55	33.64	60.26
Aya-Vision-32B	24.27	25.27	51.45	28.36	37.18	41.45	25.00	33.28
InternVL3-38B	99.00	91.36	63.45	42.64	56.73	30.73	54.64	62.65
Qwen2-VL-72B	95.55	75.36	55.27	40.64	40.55	42.36	33.55	54.75
QVQ-72B-Preview	77.82	52.18	53.09	36.73	41.82	47.73	34.64	49.14
Qwen2.5-VL-72B	99.36	87.82	65.91	45.27	43.36	58.82	41.00	63.08
InternVL3-78B	99.73	95.55	69.45	52.27	49.27	58.18	57.09	68.79

Table 2: Full Evaluation Results of 24 VLMs on Jigsaw-Puzzles. VLMs are grouped into proprietary and opensource categories. Dark Green and Light Green indicate the top-1 and top-2 performance within each group, respectively. Results of reasoning-enhanced are marked in **bold**. We also highlight the top three models based on their overall performance, using Dark Blue, Medium Blue, and Light Blue, respectively.

Modolo	Missing Pi	ece Selection	Piece Lo	calization	Connection	Anomaly	Order	Overall	
would	Easy	Hard	Easy	Hard	Verification	Detection	Restoration	Overall	
Human Performance	99.55	100.00	95.45	91.36	93.18	97.27	97.73	96.36	
Proprietary Models									
Claude-3.7-Sonnet	100.00	95.45	55.45	47.27	42.73	68.18	38.64	63.96	
Gemini-2.5-Flash	98.18	93.18	58.18	55.45	42.27	66.82	32.27	63.76	
Gemini-2.5-Flash-Thinking	99.55	95.91	71.82	51.82	55.00	60.91	57.27	70.33	
Gemini-2.5-Pro	100.00	96.36	77.73	56.82	57.27	71.36	70.91	75.78	
Open-source Models									
Qwen2.5-VL-72B	99.09	86.82	67.73	42.27	40.00	57.73	33.64	61.04	
InternVL3-78B	99.55	95.45	70.45	53.64	44.55	61.36	56.82	68.83	

Table 3: Comparing Top-Performing VLMs with Human Performance on Jigsaw-Puzzles-Lite. The human performance is highlighted in Dark Green. Results of reasoning-enhanced are marked in **bold**. The top three overall performance are highlighted in Dark Blue, Medium Blue, and Light Blue, respectively.

Proprietary VLMs. As shown in Tables 2, 370 proprietary VLMs consistently outperform open-371 source VLMs on Jigsaw-Puzzles. Among them, 372 non-reasoning-enhanced proprietary VLMs typi-373 cally exceed 60% overall accuracy, whereas most open-source VLMs fall short-only InternVL3-375 376 [14B/38B/78B] and Qwen2.5-VL-72B surpass this threshold. Reasoning-enhanced proprietary mod-377 els, such as Gemini-2.5-Flash-Thinking (72.42%) and Gemini-2.5-Pro (77.14%), further widen this 379

gap. These results reveal a persistent disparity in spatial reasoning performance, suggesting that proprietary VLMs benefit from advantages in model architecture, training strategy, and access to largescale data. Meanwhile, this finding highlights substantial room for improvement in open-source VLMs toward achieving more robust and generalizable spatial reasoning.

Model Performance in Different Tasks. In the *Missing Piece Selection* task, which primar-

381 382 383

380

388

ily targets spatial understanding, most proprietary 390 VLMs perform well under both Easy and Hard 391 settings, demonstrating strong perceptual capability. Although open source models generally perform poorly by comparison, certain models, such as the InternVL3 series and Qwen2.5-VL-72B, achieve perceptual understanding on par with pro-396 prietary VLMs. Notably, both the Aya-Vision series and the Mistral-Small-3.1-24B-Instruct models perform poorly across all settings, even at the 32B scale, accuracy remains near random, reveal-400 ing severe deficits in spatial understanding and 401 instruction following. In single-step spatial rea-402 soning tasks-Piece Localization, Connection Ver-403 ification, and Anomaly Detection-most VLMs 404 surpass the p-value-based critical value, indicat-405 ing emerging competence in basic spatial reason-406 ing. However, strong performance remains concen-407 trated in only a few models, particularly reasoning-408 enhanced proprietary models and the latest open-409 source InternVL3 series. This disparity becomes 410 even more evident in the multi-step spatial reason-411 ing task-Order Restoration, indicating that most 412 VLMs struggle with complex spatial reasoning. 413

414 In conclusion, Jigsaw-Puzzles effectively distinguishes VLMs across a spectrum of spatial reason-415 ing capability-from basic understanding to com-416 plex multi-step reasoning. As shown by the results 417 in Table 2, substantial room for improvement re-418 mains, particularly in multi-step spatial reasoning. 419 **Foundational Spatial Understanding Shapes** 420 Reasoning Performance. We analyze task similar-421 ity on Jigsaw-Puzzles by computing Pearson cor-422 relation coefficients between each task and all oth-423 424 ers, as proposed by Zhang et al. (2025), as shown in Figure 5. The results show that performance 425 on the Missing Piece Selection task-a proxy for 426 spatial understanding, is strongly correlated with 427 performance on spatial reasoning tasks. In con-428 trast, VLMs with weaker spatial understanding of-429 ten struggle with reasoning tasks, with some per-430 forming worse than random on reasoning-intensive 431 tasks. This pattern reflects the human cognitive 432 progression from perception to understanding to 433 reasoning, underscoring the foundational role of 434 spatial understanding in enabling higher-level spa-435 tial reasoning in VLMs. 436

437 Spatial Reasoning Scales with VLM size. We
438 analyze the relationship between VLM size and
439 overall performance on Jigsaw-Puzzles. As shown
440 in Figure 6, our results reveal that VLM ac441 curacy consistently increases with model size,



Figure 5: Task Similarity Heatmap. The heatmap illustrates the pairwise correlation between tasks in our benchmark, measured using Pearson correlation coefficients. Task names are abbreviated using the initials of each word (e.g., *Missing Piece Selection* \rightarrow MPS). The suffixes _e and _h indicate the Easy and Hard settings, respectively.



Figure 6: Relationship between VLM size and performance on Jigsaw-Puzzles. Each point represents a VLM, with its accuracy plotted against log-scaled parameter size. A clear positive correlation is observed both across and within model families, indicating that larger models tend to exhibit stronger performance.

both across all models and within specific families (e.g., InternVL3, Qwen2.5-VL). This positive correlation suggests that spatial reasoning capability—like other cognitive competencies (Wang et al., 2024b)—benefits from larger model capacity, which scales with parameter count.

Reasoning-Enhanced Models Show Superior Spatial Reasoning Performance. To assess the spatial reasoning capability of reasoning-enhanced VLMs, we evaluate Gemini-2.5-Flash-Thinking, Gemini-2.5-Pro, Kimi-VL-A3B-Thinking and QvQ-72B-Preview. Except for Gemini-2.5-Pro,

453



Figure 7: An example of self-correction. Red shows the initial incorrect answer generated by Gemini-2.5-Pro; Blue indicates the ground-truth answer; Green illustrates the model's self-correction process.

each model has a corresponding base version for 454 comparison. As shown in Table 2, these enhanced 455 VLMs consistently achieve higher overall accuracy. 456 For example, Kimi-VL-A3B-Thinking improves 457 from 42.04% to 44.43%, and Gemini-2.5-Flash-458 Thinking rises from 65.86% to 72.42%. Although 459 QvQ-72B-Preview overall underperforms Qwen2-460 461 VL-72B, it achieves better results on spatial reasoning tasks. Notably, Gemini-2.5-Pro achieves 462 the highest overall accuracy (77.14%) among all 463 VLMs tested. Furthermore, the largest improve-464 ments occur in the multi-step spatial reasoning 465 task, Order Restoration, where reasoning-enhanced 466 VLMs outperform their base counterparts more sub-467 stantially than in single-step tasks. To explain this, 468 we analyze cases where only Gemini-2.5-Pro an-469 swers correctly, with Figure 7 presenting one such 470 example. Gemini-2.5-Pro demonstrates a form of 471 self-correction: when the model's initial predic-472 tion is not among the provided options, it will re-473 474 evaluate the visual input and revise its judgment. This behavior, facilitated by the reduced answer 475 space under choice constraints, may contribute to 476 the superior performance of reasoning-enhanced 477 models in the Order Restoration task. 478



Figure 8: Evaluation of *Order Restoration* and *Order Generation* tasks on Jigsaw-Puzzles-Lite. Without option constraints, VLM accuracy drops significantly—peaking at just 30.00% and falling far short of human performance.

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

504

505

506

507

509

Further Exploring Multi-Step Spatial Reasoning in VLMs. To further evaluate VLMs' multistep spatial reasoning beyond the constraints of predefined choices, we introduce the Order Generation task based on Jigsaw-Puzzles-Lite. In this setting, VLMs must directly generate the correct sequence of puzzle pieces without relying on answer options, thereby more authentically simulating open-ended spatial reasoning. As shown in Figure 8, current VLMs consistently struggle with this task—Gemini-2.5-Pro, the best-performing model, achieves only 30.00% accuracy, in stark contrast to 94.09% by human participants. This finding reveals that, despite exhibiting strong self-correction behavior under option constraints, existing VLMs face considerable challenges in autonomously constructing coherent spatial reasoning chains. This highlights a significant gap between current VLMs and human-level spatial reasoning in open-ended scenarios.

5 Conclusion

We propose Jigsaw-Puzzles, a novel benchmark for systematically evaluating the spatial reasoning capability of VLMs in real-world visual scenes. Through extensive experiments on 24 representative VLMs, we identify persistent gaps between current VLMs and human-level spatial reasoning—especially in multi-step spatial reasoning tasks. Jigsaw-Puzzles provides a scalable and cognitively grounded benchmark to advance future research on spatial reasoning in VLMs.

607

608

609

610

611

612

613

614

615

616

617

510 Limitations

511While Jigsaw-Puzzles provides a structured bench-512mark tailored for 2D spatial reasoning in static im-513ages, it currently does not address 3D perception,514temporal sequences, or embodied contexts—each515of which represents an important and orthogonal516axis of spatial cognition. We view this as a natural517next step and encourage future work to extend the518benchmark in these directions.

References

520

521

522

523

524

527

528

529

530

532

533 534

535

536

538

539

540

541

543

544

545

546

549

551

552

553

554

555

557

561

- Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, and 1 others. 2025. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*.
 - Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
 - Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805.
 - Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku.
 - Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. arXiv preprint arXiv:2502.13923.
 - Nitzan Bitton-Guetta, Yonatan Bitton, Jack Hessel, Ludwig Schmidt, Yuval Elovici, Gabriel Stanovsky, and Roy Schwartz. 2023. Breaking common sense: Whoops! a vision-and-language benchmark of synthetic and compositional images. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 2616–2627.
 - Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024.
 Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465.
- Saurabh Dash, Yiyang Nan, John Dang, Arash Ahmadian, Shivalika Singh, Madeline Smith, Bharat Venkitesh, Vlad Shmyhlo, Viraat Aryabumi, Walter Beller-Morales, and 1 others. 2025. Aya vision: Advancing the frontier of multilingual multimodality. *arXiv preprint arXiv:2505.08751*.

- Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, and 1 others. 2025. Kimi-vl technical report. *arXiv preprint arXiv:2504.07491*.
- Patrick Fissler, Olivia Caroline Küster, Daria Laptinskaya, Laura Sophia Loy, Christine AF Von Arnim, and Iris-Tatjana Kolassa. 2018. Jigsaw puzzling taps multiple cognitive abilities and is a potential protective factor for cognitive aging. *Frontiers in aging neuroscience*, 10:408085.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2024. Mme: A comprehensive evaluation benchmark for multimodal large language models. *Preprint*, arXiv:2306.13394.

Grok. 2024. Bringing grok to everyone.

- Toru Ishikawa and Nora S Newcombe. 2021. Why spatial is special in education, learning, and everyday activities.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, and 1 others. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73.
- Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. 2017. Building machines that learn and think like people. *Behavioral and brain sciences*, 40:e253.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2024. Seedbench: Benchmarking multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13299–13308.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer vision– ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages 740–755. Springer.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, and 1 others. 2024. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.

- 618
- 625 627 631 632 633
- 635
- 638
- 641 642
- 643
- 645
- 647
- 648
- 654

- 667

- 671

- Mistral. 2025. Mistral small 3.1.
 - Atin Pothiraj, Elias Stengel-Eskin, Jaemin Cho, and Mohit Bansal. 2025. Capture: Evaluating spatial reasoning in vision language models via occluded object counting. arXiv preprint arXiv:2504.15485.
 - Qwen. 2024. Qvq: To see the world with wisdom.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PmLR.
 - Yufan Ren, Konstantinos Tertikas, Shalini Maiti, Junlin Han, Tong Zhang, Sabine Süsstrunk, and Filippos Kokkinos. 2025. Vgrp-bench: Visual grid reasoning puzzle benchmark for large vision-language models. arXiv preprint arXiv:2503.23064.
 - Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. Laion-400m: Open dataset of clipfiltered 400 million image-text pairs. arXiv preprint arXiv:2111.02114.
 - Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2556–2565.
 - Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruy Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vga models that can read. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 8317-8326.
- Yueqi Song, Tianyue Ou, Yibo Kong, Zecheng Li, Graham Neubig, and Xiang Yue. 2025. Visualpuzzles: Decoupling multimodal reasoning evaluation from domain knowledge. arXiv preprint arXiv:2504.10342.
- Ilias Stogiannidis, Steven McDonagh, and Sotirios A Tsaftaris. 2025. Mind the gap: Benchmarking spatial reasoning in vision-language models. arXiv preprint arXiv:2503.19707.
- Kexian Tang, Junyao Gao, Yanhong Zeng, Haodong Duan, Yanan Sun, Zhening Xing, Wenran Liu, Kaifeng Lyu, and Kai Chen. 2025. Lego-puzzles: How good are mllms at multi-step spatial reasoning? arXiv preprint arXiv:2503.19990.
- Ashish V Thapliyal, Jordi Pont-Tuset, Xi Chen, and Radu Soricut. 2022. Crossmodal-3600: A massively multilingual multimodal evaluation dataset. arXiv preprint arXiv:2205.12522.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024a. Qwen2vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191.

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

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

720

721

723

- Xinglin Wang, Peiwen Yuan, Shaoxiong Feng, Yiwei Li, Boyuan Pan, Heda Wang, Yao Hu, and Kan Li. 2024b. Coglm: Tracking cognitive development of large language models. arXiv preprint arXiv:2408.09150.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. Transactions of the association for computational linguistics, 2:67–78.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9556-9567.
- Zicheng Zhang, Xiangyu Zhao, Xinyu Fang, Chunyi Li, Xiaohong Liu, Xiongkuo Min, Haodong Duan, Kai Chen, and Guangtao Zhai. 2025. Redundancy principles for mllms benchmarks. arXiv preprint arXiv:2501.13953.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, and 1 others. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. arXiv preprint arXiv:2504.10479.

A Appendix

A.1 Examples of Other Spatial Reasoning **Benchmarks**

Figure 9 and Figure 10 illustrate two representative question types commonly used to evaluate spatial reasoning capability of VLMs. The tested images are not based on real-world scenes, which limits the capability to evaluate spatial reasoning in VLMs under realistic conditions.

A.2 Dataset Curation

Rules Pool. Figure 11 shows examples of images that were rejected and accepted based on the filtering rules, the following are the rules defined by experts during the image selection process:

- Removing images containing explicit or violent content.
- Filtering out blurry, low-resolution, or visually ambiguous images.



Figure 9: An example of Mind the Gap (Stogiannidis et al., 2025).



Figure 10: An example of LEGO-Puzzles (Tang et al., 2025).



Figure 11: Top: examples of images rejected by expertdefined filtering rules. Bottom: examples of highquality images that pass the rules.

• Excluding images lacking semantic clarity or spatial structure.

724

725

726

727

728

731

732

733

734

- Discarding images with structural ambiguity (e.g., multiple valid puzzle arrangements).
- Eliminating misaligned images or those with overly small visual elements after cropping, which hinder spatial reasoning.

Task-Specific Template. The following are detailed templates for each task. Note that <image_x> denotes a placeholder for the corresponding image input.



Figure 12: Template of *Missing Piece Selection*, notably, the templates for the Easy and Hard settings are identical.



Figure 13: Template of Piece Localization (Easy).



Figure 14: Template of Piece Localization (Hard).



Figure 15: Template of Connection Verification.

TASK 4: Anomaly Detection				
Question: You are given an image that is created by cutting a full image into 4 sub-images				
and stitching them back together. The positions of the <image 1=""/> , <image 2=""/> , <image 3=""/> ,				
<image 4=""/> in the stitched image are defined as follows:				
A. Top-left (the sub-image in the top-left corner)				
B. Top-right (the sub-image in the top-right corner)				
C. Bottom-left (the sub-image in the bottom-left corner)				
D. Bottom-right (the sub-image in the bottom-right corner)				
Note: These positions refer to the current stitched image you are looking at. Exactly one of the				
4 sub-images may have been rotated or mirrored, or all sub-images may be completely normal.				
There is at most one abnormal sub-image in each stitched image. Your tasks:				
1. Decide whether the stitched image is correct. Options: A. Correct B. Incorrect				
If incorrect, answer:				
Which sub-image is abnormal? Options: A/B/C/D				
What kind of change has occurred? Options: A. Rotation, B. Mirroring				
Output format (exactly 3 lines):				
Judgment: A or B				
Error Position: A/B/C/D (or leave blank)				
Error Type: A/B (or leave blank)				
Do not add any explanations or extra text. Example if incorrect (error in bottom-left, rotated):				
Judgment: B Error Position: C Error Type: A				
Example if correct: Judgment: A Error Position: Error Type:				

Figure 16: Template of Anomaly Detection.

TASK 5: Order Restoration
Question: You are given <image 1=""/> , <image 2=""/> , <image 3=""/> , <image 4=""/> that are cropped from an original full image. The full image was divided into four regions by splitting it through the center: top-left, top-right, bottom-left, and bottom-right. The four cropped images have been shuffled. Based on the visual content of each cropped images, determine the correct order that reconstructs the original full image. The order corresponds to the regions in the following sequence: top-left, top-right, bottom-left, hottom-right. Choose the most appropriate option
based on the mapping below:
coption list> (A. [2, 3, 1, 4] B. [2, 4, 3, 1] C. [3, 4, 2, 1] D. [4, 1, 3, 2])
Each number corresponds to the index of the shuffled images you received.
Please respond only with "A", "B", "C", or "D", without any additional explanation or

736

description

Figure 17: Template of *Order Restoration*. Note: <option list> serves as a placeholder for the answer choices. The text in parentheses is an example and should be removed in actual use. One option is the correct answer, while the remaining three are randomly drawn from the other 23 candidates.

TASK 6: Order Generation

Question: You are given <image 1>, <image 2>, <image 3>, <image 4> that are cropped from an original full image. The full image was divided into four regions by splitting it through the center: top-left, top-right, bottom-left, and bottom-right. The four cropped images have been shuffled. Based on the visual content of each cropped image, determine the correct order that reconstructs the original full image. The order corresponds to the regions in the following sequence: top-left, top-right, bottom-left, bottom-right. Please respond only with a Python-style list of four integers indicating the correct order of the shuffled images, like [2, 3, 1, 4]. Each number corresponds to the index of the shuffled images you received. Please Do not provide any explanation or description.

Figure 18: Template of Order Generation.

Hardware Setup for Evaluating VLMs. We eval-737 uate open-source VLMs using hardware configu-738 rations scaled to model size. Models with fewer than 20B parameters run on a single NVIDIA A100 740 80GB GPU. Those between 20B and 40B use two 741 NVIDIA A100 80GB GPUs, while models exceed-742 ing 40B are evaluated on four NVIDIA A100 80GB 743 GPUs to meet their greater memory and computa-744 tional demands. 745