Sparkle: Mastering Basic Spatial Capabilities IN VISION LANGUAGE MODELS ELICITS GENERAL IZATION TO COMPOSITE SPATIAL REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Vision language models (VLMs) have demonstrated impressive performance across a wide range of downstream tasks. However, their proficiency in spatial reasoning remains limited, despite its crucial role in tasks involving navigation and interaction with physical environments. Specifically, much of the spatial reasoning in these tasks occurs in two-dimensional (2D) environments, and our evaluation reveals that state-of-the-art VLMs frequently generate implausible and incorrect responses to composite spatial reasoning problems, including simple pathfinding tasks that humans can solve effortlessly at a glance. To address this, we explore an effective approach to enhance 2D spatial reasoning within VLMs by training the model on basic spatial capabilities. We begin by disentangling the key components of 2D spatial reasoning: direction comprehension, distance estimation, and localization. Our central hypothesis is that mastering these basic spatial capabilities can significantly enhance a model's performance on composite spatial tasks requiring advanced spatial understanding and combinatorial problemsolving. To investigate this hypothesis, we introduce *Sparkle*, a framework that fine-tunes VLMs on these three basic spatial capabilities by synthetic data generation and targeted supervision to form an instruction dataset for each capability. Our experiments demonstrate that VLMs fine-tuned with Sparkle achieve significant performance gains, not only in the basic tasks themselves but also in generalizing to composite and out-of-distribution spatial reasoning tasks (e.g., improving from 13.5% to 40.0% on the shortest path problem). These findings underscore the effectiveness of mastering basic spatial capabilities in enhancing composite spatial problem-solving, offering insights into systematic strategies for improving VLMs' spatial reasoning capabilities.

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

033

1 INTRODUCTION

Vision language models (VLMs) (OpenAI, 2023; Liu et al., 2023b; Chen et al., 2024c) have demonstrated near-040 human performance in tasks like image captioning (Chen 041 et al., 2015), visual question answering (VQA) (Goyal 042 et al., 2017; Singh et al., 2019) and abundant down-043 stream tasks by combining visual and text inputs to reason 044 about the physical world. However, these models exhibit significant limitations in understanding spatial relationships. For instance, as shown in Figure 1, state-of-the-046 art (SoTA) VLMs GPT-40 and InternVL2-Pro (OpenAI, 047 2023; Chen et al., 2024c) generate implausible responses 048 to a shortest path problem that a human could solve at a glance, a simple 2D spatial reasoning task. 050

Nevertheless, 2D spatial reasoning is essential for VLMs
to understand and interact with the physical environments, shaping their ability to solve mazes (Ivanitskiy et al., 2023; Wang et al., 2024), plan routes (Feng et al.,



Figure 1: SoTA VLMs fail to solve the pathfinding problem, a simple 2D spatial reasoning task.



Figure 2: An overview of the workflow of the Sparkle framework.

2024; Chen et al., 2024b), and solve geometric problems like humans (Fernandes & de Oliveira, 2009). These tasks emphasize 2D spatial reasoning, requiring VLMs to process and navigate flat visual planes, interpret spatial relationships, and make decisions based on geometric understanding.
Such capabilities are fundamental in translating visual input into actionable insights. While more and more VLMs are developed with larger training datasets and extensive benchmarks (Ge et al., 2024; Zhang et al., 2024), the focus on enhancing spatial reasoning has received comparatively less attention, despite its importance to the core capabilities of VLMs.

In this paper, we study VLMs' spatial reasoning capabilities in a 2D space by investigating three key questions: (1) How well do existing models perform on 2D spatial reasoning? (2) What fundamental tasks affect spatial reasoning capabilities in 2D? (3) Can mastering basic tasks help improve the performance of complex spatial reasoning?

076 We first analyzed various spatial reasoning tasks presented in existing works (Kamann & Rother, 077 2020), identifying the capabilities required for these tasks. From this analysis, we identified three basic capabilities fundamental for spatial reasoning in 2D space: direction comprehension, distance 079 estimation, and localization. A systematic evaluation of the performance of existing open-source and closed-source VLMs on these three basic capabilities reveals that even the most advanced VLMs 081 sometimes struggle with these fundamental tasks. For instance, in a simple 2D direction classification task, where a model is asked to determine the relative direction (top left, top right, bottom left, 083 bottom right) of one object relative to another on a straightforward diagram with only two objects, the state-of-the-art VLM GPT-40 can achieve only 76.5% accuracy. In contrast, a human should be 084 able to answer these questions correctly without much thought. 085

Most real-world spatial reasoning tasks, such as pathfinding (Lester, 2005; Cui & Shi, 2011), inher ently require the composition of the basic capabilities identified above. A composite task is often
 subject to specific constraints that necessitate tailored solutions, unlike improving basic spatial reasoning capabilities, which can exhibit generalizability. In order to effectively improve the model's
 overall spatial reasoning capabilities in 2D space, we raise a conjecture: whether a VLM that masters the three basic capabilities can generalize and perform better on more complex composite spatial
 tasks. In other words, can a VLM exhibit compositional generalizability (van Zee, 2020) in spatial
 reasoning tasks?

To test this, we propose *Sparkle*, as shown in Figure 2, a framework that fine-tunes VLMs on these three basic spatial capabilities by programmatically generating synthetic data and providing supervision to form an instruction dataset for each capability. Our experimental results show that models trained on *Sparkle* achieve significant performance gains, not only in the basic tasks themselves (e.g., improving from 35% to 83% for InternVL2-8B on direction comprehension) but also in generalizing to composite and out-of-distribution general spatial reasoning tasks (e.g., improving from 13.5% to 40.0% on the shortest path problem). Additionally, our ablation study confirms the importance of mastering all three basic spatial reasoning capabilities. To summarize, our contributions are:

- We show that state-of-the-art VLMs struggle with composite spatial reasoning tasks that humans solve effortlessly.
- We identify three key components of spatial reasoning and construct an instruction-tuning method called Sparkle to improve these three fundamental spatial reasoning capabilities.
- Our experiments demonstrate that enhancing VLMs' basic spatial capabilities significantly improves their ability to generalize to out-of-distribution composite spatial tasks.

108 2 RELATED WORK

110 2.1 VISION LANGUAGE MODELS AND APPLICATIONS

112 Early works on VLMs, such as CLIP (Radford et al., 2021a) and ALIGN (Jia et al., 2021), leveraged 113 contrastive learning to align visual and textual embeddings in a shared latent space, demonstrating 114 strong capabilities in linking visual content with corresponding natural language descriptions. With the rapid advancement of Large Language Models (LLMs), modern VLMs increasingly combine 115 116 pretrained vision models (Dosovitskiy et al., 2021; Chen et al., 2023) with powerful LLMs (Chiang et al., 2023; Bai et al., 2023a; Jiang et al., 2023; Cai et al., 2024) to facilitate a more cohesive under-117 standing of both modalities (Liu et al., 2023b; Bai et al., 2023a; Chen et al., 2024c). This approach 118 enables richer visual reasoning, open-ended image captioning, and more interactive multimodal di-119 alogue systems. 120

VLMs have been applied in various pre-training tasks, such as image-text matching, masked image
modeling, and multimodal reasoning (Li et al., 2022; 2023a; Wang et al., 2022b). In downstream
tasks, they excel in applications like visual question answering (Antol et al., 2015; Wang et al., 2022a), image captioning (Li et al., 2020; Sidorov et al., 2020; Wang et al., 2021), image generation
based on textual prompts (Ramesh et al., 2022; Baldridge et al., 2024), and aiding human-machine
interactions in complex real-world settings, showcasing their versatility and potential across a broad
range of vision language applications.

- 128
- 129 2.2 SPATIAL REASONING IN LLMS AND VLMS

130 Spatial reasoning in LLMs involves understanding and manipulating spatial relationships described 131 in text. Early work focused on extracting spatial information from natural language (Hois & Kutz, 132 2011; Kordjamshidi et al., 2011). More recent efforts emphasize improving multi-hop spatial rea-133 soning (Li et al., 2024b), especially in complex scenarios like 2D visual scenes (Shi et al., 2022). 134 Key methods include pretraining on synthetic datasets to better capture spatial patterns (Mirzaee 135 et al., 2021), and using in-context learning to generalize spatial reasoning across tasks, such as 136 transforming spatial data into logical forms or visualizing reasoning traces (Yang et al., 2023b; Wu 137 et al., 2024; Tang et al., 2024).

138 Building on these foundations, VLMs extend spatial reasoning by integrating visual inputs and often 139 implicitly encode spatial knowledge through large-scale pretraining on visual-text datasets (Radford 140 et al., 2021b; Li et al., 2023b). Early studies on VLMs primarily focus on understanding spatial 141 relationships between objects in front-view images (Liu et al., 2023a), laying the groundwork for 142 2D spatial reasoning. More recently, research on VLMs has expanded to 3D reasoning tasks, which 143 introduce additional challenges such as depth estimation (Chen et al., 2024a) and path planning (Chen et al., 2024b; Deng et al., 2020), as seen in applications like robotic grasping (Xu et al., 2023) 144 and navigation (Shah et al., 2023; Chiang et al., 2024) in the embodied AI field (Li et al., 2024c). 145 Despite these advances, 2D spatial reasoning remains more fundamental and flexible, as it can be 146 applied to various tasks, including VQA (Ge et al., 2024; Kamath et al., 2023; Li et al., 2024a) and 147 user interface grounding (Rozanova et al., 2021). Due to its broad applicability and foundational 148 role, this work focuses on exploring 2D spatial reasoning capabilities within VLMs. 149

150 151

3 Methodology

152 153

154

155

156

157

In order to systematically evaluate and enhance the spatial reasoning capabilities of VLMs in 2D environments, we introduce the Sparkle framework, as illustrated in Figure 3. This section is structured as follows:

- Disentangling basic elements: How we identified the basic spatial capabilities of 2D spatial reasoning, and why these elements are foundational.
- Sparkle: We present the Sparkle framework to enhance VLMs' performances in 2D spatial reasoning by systematically improving the identified basic spatial capabilities.
- Tasks: We employ three spatial reasoning tasks specifically designed to evaluate both basic and composite spatial reasoning capabilities of VLMs.

164	Task	Recognition	Counting	Depth	Direction	Localization	Distance
165	QualSR (Freksa, 1991)	V	×	×	v	~	~
166	MapVQA (Wang et al., 2024)	~	v	×	v	V	×
167	NavVQA (Wang et al., 2024)	~	~	×	~	 Image: A set of the set of the	×
107	GridVQA (Wang et al., 2024)	~	~	×	~	 ✓ 	×
168	VisualSR (Rajabi & Kosecka, 2023)	~	×	~	~	×	~
169	TvRecog. (Li et al., 2024a)	~	×	×	×	×	×
170	TvLoc. (Li et al., 2024a)	~	×	×	×	 Image: A start of the start of	×
	StaticSR (Li et al., 2024a)	~	~	×	~	×	×
171	DynamicSR (Li et al., 2024a)	~	~	×	~	×	×
172	SRR (Chen et al., 2022)	~	×	~	×	~	V
173	COCO-Spatial (Ranasinghe et al., 2024)	~	×	 Image: A start of the start of	~	 Image: A start of the start of	×
	What's Up (Kamath et al., 2023)	~	×	 ✓ 	~	×	×
174	Q-Spatial (Liao et al., 2024)	~	×	~	×	×	~
175	SpatialRGPT (Cheng et al., 2024)	~	×	~	v	~	~

Table	1:	Overview	of VI	LM	tasks	related	to	spatial	reasoning	and	their rec	juired ca	apabilities.	•
													1	

162 163

3.1 DISENTANGLING SPATIAL REASONING

179 To systematically disentangle basic 2D spatial reasoning capabilities, we first analyze the capabilities required in existing VLM benchmarks related to spatial reasoning, as shown in Table 1. While 181 these benchmarks include a wide array of capabilities, such as image recognition and depth esti-182 mation, we narrow our focus to those most fundamental to 2D spatial reasoning. Depth estimation, 183 though relevant to spatial reasoning, is more suited to 3D tasks and thus excluded from our analysis, 184 as discussed in Section 2.2. We present the definitions of three basic components: (1) Direction 185 *Comprehension*: The ability to understand the orientation of an object relative to a reference object; (2) *Distance Estimation*: The ability to gauge the magnitude of spatial displacement between objects; (3) *Localization*: The ability to determine the precise position of an object in space. 187

188 The selected basic spatial reasoning capabilities are foundational because they collectively represent 189 the minimal components necessary to fully describe an object's position in 2D space. In particular, 190 each of the three capabilities aligns with principles from Cartesian and polar coordinate systems, 191 which serve as the mathematical bedrock for spatial representation: direction defines orientation, 192 distance represents magnitude, and localization integrates both to precisely define an object's position, ensuring comprehensive spatial awareness (Zeng & Si, 2017). This decomposition enables 193 a systematic evaluation of spatial reasoning by isolating the key dimensions of spatial understanding. 194

195 196 197

3.2 Sparkle

To comprehensively investigate our hypothesis, we introduce Sparkle, a simple yet effective frame-199 work for constructing an instruction dataset focused on enhancing a model's spatial reasoning abili-200 ties. This framework only improves VLMs' basic spatial capabilities, and this design enables us to 201 evaluate whether models that perform well on basic spatial reasoning tasks can also excel in more 202 complex and composite problems.

203

204 INSTRUCTION DATA GENERATION 3.2.1 205

206 The design of our instruction dataset focuses on three basic spatial capabilities: direction, distance, and localization, based on insights provided in Section 3.1. The proposed fine-tuning pipeline does 207 not require manual labeling, as all data can be programmatically generated. 208

209 We use \mathbb{G} to denote a data generator that can generate a set of objects, $P = \{N_i\}_{i=1}^n$, representing a training sample of basic spatial capabilities. Each object $N_i = (x_i, y_i) \in \mathbb{R}^2$ consists of randomly 210 sampled coordinates within a bounded region. For each basic capability $\mathcal{T} \in \{\text{dir., dist., loc.}\}$, we 211 construct a dataset $D_{\mathcal{T}}$ containing input-output pairs $(\mathcal{X}^{\mathcal{T}}, \mathcal{Y}^{\mathcal{T}})$, where $\mathcal{X}^{\mathcal{T}}$ represents the inputs 212 and $\mathcal{Y}^{\mathcal{T}}$ represents the corresponding ground truth outputs. Each input $\mathcal{X}^{\mathcal{T}}$ consists of: (1) A visual 213 input $\mathcal{X}_V^{\mathcal{T}}$: A labeled diagram representing the spatial configuration of a sample of objects through 214 a visual representation function $\mathbb{V}_{\mathcal{T}}(P)$, (2) A language prompt $\mathcal{X}_{L}^{\mathcal{T}}$: A question querying some 215 aspects of the spatial properties for P.

217

218

219

220

221

222

224

225 226

227

228

229

230

231

233

234

235

237 238

239 240 241

242

243

244

245

246

253

254 255

256

257

258 259 260

261

262

263 264 265

266

Sparkle Instruct Basic Spatial Capabilities Raw Data Composite NI Spatial Reasoning NB Image: 10×10 Shortest Path \$ Improve Proble $\mathbb{V}_{\mathcal{T}}$ Object Coordinates: N1: (7.12, 9.35) N2: (6.46, 2.08) Traveling N2 N3: (3.59, 7.34) Salesman Problem Basic Direction Distance Localization Figure 3: The proposed Sparkle framework. **Visual Representation** Direction Localization Distance 0: Which relative location the N2 located at? Determine the Q: Which distance is the Q: direction from the N2 object to the N3 object. shortest? N1 and N2, N1 and N3, N2 and N3 A: down right N1 A: top left Q: Which relative position A: N1 and N3 is the N1 located at? N3 Q: From the N1 object to A: top right Compare the distances: the N2 object, which N1 and N3 and N2 and N3. Q: Identify the location direction should you move Which one is longer? the N3. A: down of A: N2 and N3 A: top left Q: What is the direction Q: In a 10x10 image, from the N1 object to the Q: In a 10x10 image, what what is the **distan** N2 N3 object? is the coordinate of the between the N1 and N3 A: down left N3 object? objects? A: (3.59, 7.34) A: 4.07

Figure 4: A data sample from the Sparkle dataset.

For example, to craft a training sample for direction comprehension, two objects, N_1 and N_2 , are selected from P, and a question such as "What is the direction of N_2 relative to N_1 ?" is posed. The corresponding correct answer $Y^{\mathcal{T}}$ can be easily computed since we can access the exact coordinates of these objects, e.g., we can obtain the answer to the above question by calculating the vector from N_1 to N_2 based on their coordinates and map it to the corresponding directional label. Details about these generation processes can be found in Appendix §A.1.

247 The resulting training dataset consists of these generated questions and answers, paired with the 248 corresponding visual representations, as shown in Figure 4. Specifically, the training pairs are rep-249 resented as $\{(\mathcal{X}_L^{\text{train}}, \mathcal{X}_V^{\text{train}}, \mathcal{Y}^{\text{train}})\}$, where $\mathcal{X}_L^{\text{train}}$ represents the language-based queries, $\mathcal{X}_V^{\text{train}}$ repre-250 sents the visual representations, and \mathcal{Y}^{train} represents the corresponding answers. We also provide a 251 complete data sample from the Sparkle training set in Appendix §A.5.1. 252

3.2.2 INSTRUCTION FINETUNING FOR BASIC TASKS

To enhance the spatial reasoning capabilities of VLMs, we use the Sparkle training set, denoted as $\mathcal{X}^{\text{train}} = \{(\mathcal{X}_{L}^{\text{train}}, \mathcal{X}_{V}^{\text{train}})\}$. The objective is to minimize the negative log-likelihood of the predicted answers. Specifically, the loss function \mathcal{L} is defined as:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(\mathcal{X}^{\text{train}}, \mathcal{Y}^{\text{train}})} \left[\log p(\mathcal{Y}^{\text{train}} \mid \mathcal{X}_{V}^{\text{train}}, \mathcal{X}_{L}^{\text{train}}; \theta) \right]$$

where θ represents the parameters of the VLM. The training aims to improve the model's proficiency in basic spatial reasoning tasks, which subsequently allows for an effective evaluation of its performance on more complex spatial challenges.

3.3 TASKS

The goal of the employed tasks is to evaluate the 2D spatial reasoning capabilities of VLMs and 267 provide a foundation for studying how acquiring basic spatial capabilities can enhance performance 268 on complex tasks. To achieve this, we follow key design criteria: (1) focus on spatial reasoning, and 269 (2) progression from basic to composite tasks.

287

289

290

299

300

306



Figure 5: A sample from the evaluation dataset.

3.3.1 BASIC TASKS: ASSESSING FUNDAMENTAL SPATIAL CAPABILITIES

As shown in Figure 5 (left), the basic tasks in Sparkle are designed to assess the model's understanding of three basic spatial capabilities: (1) direction comprehension, (2) distance estimation, (3) localization.

In each basic task, the VLM is provided with an image containing several labeled data objects and a multiple-choice question about the spatial properties of these objects, with the goal of having the model answer these questions correctly. We first generate labeled diagrams that serve as visual inputs, then generate the questions (in multiple-choice format) and corresponding answer pairs, similar to the process described in Section 3.2, to obtain the basic task test set.

3.3.2 COMPOSITE TASKS: EVALUATING INTEGRATED SPATIAL REASONING

Building on the basic spatial relationships, the composite tasks introduce greater complexity. The objective here is to assess whether the model can apply basic spatial skills to solve problems requiring a combination of these skills or whether it has merely learned each skill in isolation without being able to generalize effectively. We choose the *Shortest Path Problem (SPP)* and the *Traveling Salesman Problem (TSP)* as composite tasks to evaluate the integration of basic capabilities.

Shortest Path Problem (SPP) As shown in Figure 5 (middle), SPP evaluates the model's capability to compute the most efficient route between two objects on a 2D grid, requiring a combination of distance estimation and spatial planning.

Consider a grid G of size $n \times n$, with two special objects: the start object $N_{\text{start}} = (x_s, y_s)$ and the end object $N_{\text{end}} = (x_e, y_e)$. We employ a language model LM generates the prompt $\mathcal{X}_L^{\text{spp}}$ using a predefined prompt template \mathbb{P}_{spp} , expressed as: $\mathcal{X}_L^{\text{spp}} = \text{LM}(\mathbb{P}_{\text{spp}}(G, N_{\text{start}}, N_{\text{end}}))$. The visual input is produced similar to basic tasks: $\mathcal{X}_V^{\text{spp}} = \mathbb{V}_{\text{spp}}(G, N_{\text{start}}, N_{\text{end}})$.

The combined input for the VLM is $\mathcal{X}^{\text{spp}} = (\mathcal{X}_V^{\text{spp}}, \mathcal{X}_L^{\text{spp}})$, and the model is expected to predict the shortest path $\hat{\mathcal{Y}}^{\text{spp}}$, which is evaluated against the true shortest path, \mathcal{Y}^{spp} , computed using standard algorithms.

- Traveling Salesman Problem (TSP) As shown in Figure 5 (right), the TSP represents a more challenging spatial reasoning task, involving combinatorial optimization. The model must find the shortest possible route that visits each object exactly once and returns to the starting object.
- Given *n* objects $P^{\text{tsp}} = \{N_i\}_{i=1}^n$ sampled from \mathbb{G} , the ground truth solution \mathcal{Y}^{tsp} is computed using a TSP solver $\mathbb{M}_{\text{tsp}}(P^{\text{tsp}})$. Similarly, the input to VLMs consists of a visual representation $\mathcal{X}_V^{\text{tsp}} = \mathbb{V}_{\text{tsp}}(P^{\text{tsp}})$ and a corresponding language prompt $\mathcal{X}_L^{\text{tsp}}$. The complete input query is $\mathcal{X}^{\text{tsp}} =$

 $\begin{array}{l} \overset{324}{} & (\mathcal{X}_V^{\mathrm{tsp}}, \mathcal{X}_L^{\mathrm{tsp}}). \text{ Similarly, the model's predicted order of visiting all objects, } \widehat{\mathcal{Y}}^{\mathrm{tsp}}, \text{ is then evaluated} \\ \overset{325}{} & \text{against the ground truth solution } \mathcal{Y}^{\mathrm{tsp}}. \end{array}$

3.3.3 DISCUSSION

Given that the SPP can be solved in polynomial time, we expect that if the model can effectively combine its knowledge of basic spatial concepts, it will show significant improvements in solving this task efficiently. On the other hand, the TSP is an NP-hard problem, requiring combinatorial optimization to obtain the exact solution. We include the TSP to push the limits of the model's spatial reasoning capabilities, aiming to investigate how well the model can manage more complex problem-solving tasks beyond the basic integration of spatial skills.

335 336

337

344

327

328

4 EXPERIMENTS

In this section, we provide our findings along with the supporting results to demonstrate the effectiveness of the Sparkle framework. Specifically, the experiments are designed to answer the following research questions: **RQ1**: Can mastering basic 2D spatial components enhance overall spatial reasoning capability in VLMs? **RQ2**: What insights from the results of evaluations (Section 4.2.1), enhancements (Section 4.2.2), and spatial components (Section 4.3) can guide improvements in model design, training strategies, and data collection for spatial reasoning in VLMs?

344 345 4.1 SETTINGS

346 **Models** We tested open-source and commercial models to evaluate and enhance VLMs' spatial 347 reasoning capabilities. For commercial VLMs, we used GPT-40 from OpenAI (Yang et al., 2023a) 348 and Google-Gemini (GeminiTeam et al., 2023). We included LLaVA1.6 (Liu et al., 2024) and 349 InternVL2 (Chen et al., 2024c) for open-source models. Detailed model specifications and configu-350 rations are provided in Appendix §A.1. We use the MS-Swift library (Zhao et al., 2024) and apply 351 the LoRA (Hu et al., 2022) fine-tuning strategy, with low-rank dimension of 32. We set a constant 352 learning rate of 1e-4 and a batch size of 1. All training and evaluation tasks are performed on GPU 353 clusters with $8 \times NVIDIA$ A100 machines. Further details can be found in Appendix §A.1.

354

355 **Data** We built the Sparkle training dataset by generating 2,000 images. each with 17 instructionanswer pairs that describe the spatial relationships between objects, resulting in 34K samples in 356 total. Out of these 17 pairs, 3 focus on directions between objects, 7 on distances (including 4 for 357 comparing distances and 3 for estimating numerical distances), and 6 relate to localization (with 3 for 358 identifying object locations and 3 for estimating exact positions). The final instruction describes the 359 overall spatial relationships in the image. This setup helps ensure the VLM maintains its capability 360 to follow instructions effectively. Our evaluation includes tasks of: (1) shortest path problem (SPP), 361 (2) traveling salesman problem (TSP), and (3) basic spatial relationship understanding. For each 362 of them, we generated 200 samples, which together make up the evaluation set. For SPP and TSP, 363 we use LLaMA 3.1 (Dubey et al., 2024) to process the VLMs' responses into list formats to enable 364 metric computation. For the basic tasks, we structured them as a multiple-choice question format. 365 In addition, for SPP and TSP, we designed experiments that vary by grid size and the number of 366 objects involved. Detailed data statistics and sample data are provided in Appendix §A.2.

To further assess the generalizability of the improved spatial reasoning capabilities, we evaluated VLMs on existing general spatial reasoning-related benchmarks to examine its out-of-distribution performance. The general spatial benchmarks we used include What's Up, COCO-spatial, and GQA-spatial (Kamath et al., 2023), which feature real-world images and spatial reasoning questions.

371

Metrics For most of the tasks, we report accuracy as the primary evaluation metric. However, for tasks like SPP and TSP, where distance is crucial, we also use an additional metric *Margin* as a complementing one. This metric measures the extent to which the total distance of the solved path exceeds the optimal path (i.e., the shortest), expressed as a ratio of their summed distances. A lower Margin indicates better performance. Formally, the Margin is defined as: Margin = $\sum (dist(solved_path) - dist(optimal_path)) / \sum dist(optimal_path)$, where the function dist(·) computes the total distance of a given path. Table 2: VLMs' performance on basic and composite spatial reasoning tasks. "SPP-*n*Grid" denotes shortest path problem on $n \times n$ grids. "TSP-*n*Obj" denotes traveling salesman problem with *n* objects.

Model	Basic Tasks			SPP-4Grid		SPP-5Grid		TSP-4Obj		TSP-5Obj	
Woder	Loc.	Dist.	Dir.	Acc↑	$Margin \downarrow$	Acc↑	Margin↓	Acc↑	Margin↓	Acc↑	Margin↓
GPT-40	67.5	41.5	76.5	74.5	0.089	78.5	0.178	23.5	0.001	21.5	0.195
Gemini	61.5	40.5	55.0	67.0	0.208	65.0	0.188	11.5	0.023	21.7	0.070
LLaVA1.6-7B	24.5	37.0	30.5	1.5	-	0.0	-	16.0	0.132	5.0	0.476
InternVL2-26B	62.5	45.5	58.0	15.5	0.5	10.9	1.135	21.5	0.074	12.5	0.265
InternVL2-8B	60.5	44.5	35.0	16.5	3.893	13.5	1.127	17.5	0.073	11.5	0.302
+ Sparkle-Instruct	73.0	84.0	83.0	36.5	0.571	40.0	0.466	20.0	0.043	14.5	0.239
Δ	+21%	+89%	+137%	+121%	-85%	+196%	-58%	+14%	-41%	+26%	-21%

390 391 392

393

394

381 382

4.2 MAIN RESULTS

4.2.1 EVALUATION OF EXISTING VLMS

395 From Table 2, we observe that even the state-of-the-art commercial VLMs cannot obtain satisfactory 396 results on composite tasks like SPP and TSP. Open-source models achieve even worse performance 397 $(\leq 25\%$ accuracy) on these tasks. Specifically, LLaVA perform poorly particularly on SPP com-398 pared to TSP, which may attribute to the grid data structure in SPP is more complex for VLMs to 399 perceive compared to handling just a few objects in TSP, indicating that these VLMs struggle with 400 visual representations involving intricate spatial structures. Performance on the TSP task worsens 401 as the number of objects increases across most models, highlighting the growing difficulty of spatial reasoning with more objects. However, in SPP, increasing the grid size has little impact on perfor-402 mance, indicating that a larger grid does not increase the difficulty of reasoning. This result aligns 403 with our initial design principles, where SPP was intended to combine basic spatial relationship 404 understanding with a straightforward form of spatial planning. 405

To delve into how VLMs behave poorly on the composite spatial reasoning tasks, we further examine their performance on basic spatial relationship understanding, i.e. direction, location and localization comprehension. As shown in Table 2, even the state-of-the-art VLM GPT-40 struggles with basic spatial relationship understanding, achieving only 67.5%, 41.5%, and 76.5% accuracy on the direction, distance, and localization tasks, respectively. This investigation helps explain why VLMs underperform on composite tasks, as their inadequate basic spatial reasoning capabilities directly hinder their ability to handle more complex spatial challenges.

- 413
- 414 4.2.2 EFFECTIVENESS OF SPARKLE

To demonstrate the effectiveness of Sparkle, we present results from fine-tuning InternVL2-8B with this method. The results reveal significant improvements in both basic and composite tasks, indicating that 2D spatial reasoning capabilities can be significantly improved when a model effectively masters the basic components of 2D spatial reasoning.

Specifically, Sparkle only contains data for basic spatial relationship understanding. However, after 420 fine-tuning with this data, VLMs improved in basic spatial reasoning (around 80%) and showed 421 significant gains (around 90%) in composite tasks. This justifies the soundness of our abstraction 422 of spatial reasoning in 2D space into three basic components (i.e., localization, distance, and direc-423 tion), and that improving these basic capabilities could effectively enhance VLMs' overall spatial 424 reasoning, enabling it to tackle more complex tasks. This finding highlights the potential of strength-425 ening basic capabilities to improve problem-solving performance in VLMs. When comparing the 426 improvements of the InternVL2 model on SPP and TSP, we observe that the gains (around 20%) on 427 TSP are much smaller than those on the SPP task (160%). One possible explanation is that the TSP 428 involves more complex optimization challenges, which may not be as easily addressed by simply 429 improving basic spatial reasoning skills, as discussed in Section 3.3.3. This underscores the need for further research into the optimization capabilities of language models, a topic we hope our find-430 ings will inspire. Additionally, we present results of Sparkle on Qwen-VL-7B in Appendix §A.3.1, 431 demonstrating its effectiveness across various different VLMs.

Model	What's Up	COCO-	Spatial	GQA-	Spatial	
Model		1Obj	2Obj	10bj	2Obj	
GPT-40	95.9	88.2	49.7	89.4	63.6	
Gemini	69.4	50.8	34.1	42.9	21.7	
LLaVA1.6-7B	44.9	14.4	6.0	12.6	2.2	
InternVL2-26B	87.9	72.7	62.9	91.4	75.0	
InternVL2-8B	92.7	92.5	71.3	97.5	85.3	
+ Sparkle-Instruct	93.9	93.0	78.4	98.0	90.0	
Δ	+1.3%	+0.5%	+10%	+0.5%	+5.5%	

Table 3: Results on general spatial tasks.

4.2.3 GENERALIZABILITY

In the previous subsection, we have shown that spatial reasoning improvements can generalize from
 simple tasks to more complex ones. In this section, we evaluate this generalization further by testing
 spatial reasoning performance in an out-of-distribution (OOD) visual representation setting.

446 Specifically, we investigate whether the enhanced spatial reasoning capabilities transfer to other 447 general VLM spatial tasks. As seen in Table 3, there are consistent gains across general VLM 448 benchmarks related to spatial reasoning. For instance, the COCO-spatial and GQA-spatial bench-449 marks illustrate that current VLMs often struggle to accurately capture spatial relationships between 450 two objects. However, with our Sparkle framework, this capability is greatly improved.

These findings suggest that future work designing and training VLMs should consider improving
 spatial reasoning of VLMs by decomposing into basic capabilities to enhance the general performance. Our results demonstrate that the Sparkle framework is simple and highly effective in enhancing spatial reasoning capabilities in VLMs.

- 4.3 ABLATION STUDIES
- 456 457 458 459

460 461

462

455

441

442

In this section, we present the results of ablation studies on the proposed Sparkle framework, using the InternVL2-8B model for demonstration.

4.3.1 IMPACT OF TRAINING COMPONENTS

To evaluate the impact of different training 463 components, we compared Sparkle to sev-464 eral variants. First, we trained InternVL2-465 8B on individual spatial reasoning tasks 466 with our Sparkle framework, resulting 467 in Sparkle(Direction, Distance, Localiza-468 tion). Additionally, we tested a version 469 called Sparkle w/o Num that excludes nu-470 merical information (i.e., distance and lo-471 cation estimation) in Sparkle. All of the four variants are trained with the same 472 number of total samples as the full Sparkle 473 model. The results shown in Figure 6 474 reveal two key insights: First, Sparkle 475 w/o Num consistently underperforms com-476 pared to the full Sparkle model, particu-477 larly in tasks that require strong distance 478 reasoning, such as TSP. This suggests that 479 incorporating numerical information dur-



Figure 6: Ablation results showing accuracy for different Sparkle variants: Sparkle ; Sparkle without numerical information ; Sparkle (Localization) ; Sparkle (Distance) ; Sparkle (Direction) .

ing training significantly enhances the model's capability in tasks involving distance reasoning and
other related composite challenges. Second, training on specific spatial reasoning subsets can sometimes yield optimal performance for certain tasks. For example, *Sparkle (Direction)* achieves 96.4%
accuracy on the What's Up benchmark, indicating that task-specific training can be highly effective. This highlights the importance of tailoring the training process to the unique characteristics
of individual tasks. When a task emphasizes a particular spatial reasoning capability, focusing the
training data on that aspect can improve performance on the targeted task. Overall, the full Sparkle



Table 4: Results of Sparkle on InternVL2-8B with varying training sample sizes.

framework consistently delivers the best results across the majority of benchmarks, demonstrating the effectiveness of a more comprehensive approach to training.

4.3.2 IMPACT OF TRAINING SAMPLE SIZE

500 In this section, we varied the training sample size in Sparkle and evaluated its impact on spatial 501 reasoning tasks. The results are shown in Figure 4. Several key trends emerge from the results. 502 First, we observe a general improvement in VLM performance as the training sample size increases 503 despite some fluctuations in the curve. However, a noteworthy finding is the existence of task-504 specific sweet spots, beyond which performance gains taper off or degrade. This suggests that scaling up training samples does not always yield proportional improvements. For example, in the 505 TSP task, performance begins to degrade once the number of training samples surpasses a threshold 506 (around 800). This is likely because, after mastering basic spatial relationships, the model may 507 focus on locally optimal choices, such as selecting the nearest objects to form a path, rather than 508 optimizing the entire path. As the model grows more confident in these local decisions, it may 509 sacrifice the global optimality of the solution, resulting in suboptimal performance. 510

511 4.4 DISCUSSION 512

513 The analysis and results confirm that mastering basic 2D spatial reasoning capabilities through 514 Sparkle can significantly enhance VLMs' overall spatial reasoning in composite tasks (e.g., spa-515 tial planning) and general spatial tasks. This directly addresses RQ1 and supports the assumption 516 presented in the methods section.

517 Turning to RQ2, the evaluation results revealed the limitations of existing VLMs, particularly in their 518 capability to perceive complex spatial structures, as evidenced in tasks like SPP. This highlights the 519 need for improved model and training designs to support more detailed spatial reasoning. Moreover, 520 introducing synthetic data focusing on basic spatial relationships has proven to enhance overall 521 VLM spatial reasoning performance, offering a clear path for future spatial data collection. Lastly, 522 our ablation study suggests that training specific spatial reasoning capabilities in isolation yields the 523 best results for tasks that demand focused spatial abilities. Therefore, in terms of training strategy, our findings suggest adopting a pre-train and fine-tune approach (i.e., using diverse spatial data 524 in pretraining and fine-tuning specific spatial capabilities tailored to particular tasks) to improve 525 VLMs' performances on corresponding tasks. 526

527

494 495

496

497 498

499

- 5 CONCLUSION
- 528 529

530 This work presents the Sparkle framework to address the relatively limited spatial reasoning ability 531 of Vision Language Models (VLMs). Sparkle is designed to enhance spatial reasoning by focusing on three fundamental capabilities: direction comprehension, distance estimation, and localiza-532 tion. Our experiments demonstrate that fine-tuning on these basic capabilities leads to substantial 533 improvements not only in the basic tasks but also in more complex, composite spatial reasoning 534 challenges, thereby showcasing the compositional generalizability of our method. Furthermore, our 535 analysis confirms that mastering all three basic spatial reasoning capabilities is essential for broader 536 generalization, ultimately strengthening VLMs' ability to interact with the physical world. 537

- 538

540 REFERENCES

580

581

582 583

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zit nick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pp. 2425–2433, 2015.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. arXiv: 2309.16609, 2023a.

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023b.
- Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, Kelvin Chan, Yichang Chen, Sander Dieleman, Yuqing Du, Zach Eaton-Rosen, et al. Imagen 3. arXiv preprint arXiv:2408.07009, 2024.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui 561 Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye 562 Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting 563 Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, 564 Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun 565 Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang 566 Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, 567 Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, 568 Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, and 569 et al. InternIm2 technical report. arXiv: 2403.17297, 2024. 570
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia.
 Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14455–14465, 2024a.
- Jiaqi Chen, Bingqian Lin, Ran Xu, Zhenhua Chai, Xiaodan Liang, and Kwan-Yee Wong. Mapgpt: Map-guided prompting with adaptive path planning for vision-and-language navigation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9796–9810, 2024b.
 - Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. Language conditioned spatial relation reasoning for 3d object grounding. *Advances in neural information processing systems*, 35:20522–20535, 2022.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv* preprint arXiv:1504.00325, 2015.
- ⁵⁸⁷ Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv:* 2312.14238, 2023.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong,
 Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv:2404.16821*, 2024c.

612

626

632

639

640

- An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision language model. *arXiv preprint arXiv:2406.01584*, 2024.
- Hao-Tien Lewis Chiang, Zhuo Xu, Zipeng Fu, Mithun George Jacob, Tingnan Zhang, TsangWei Edward Lee, Wenhao Yu, Connor Schenck, David Rendleman, Dhruv Shah, et al. Mobility
 vla: Multimodal instruction navigation with long-context vlms and topological graphs. *arXiv preprint arXiv:2407.07775*, 2024.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https: //lmsys.org/blog/2023-03-30-vicuna/.
- Kiao Cui and Hao Shi. A*-based pathfinding in modern computer games. *International Journal of Computer Science and Network Security*, 11(1):125–130, 2011.
- Zhiwei Deng, Karthik Narasimhan, and Olga Russakovsky. Evolving graphical planner: Contextual
 global planning for vision-and-language navigation. *Advances in Neural Information Processing Systems*, 33:20660–20672, 2020.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at
 scale. In *ICLR*, 2021.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Jie Feng, Yuwei Du, Tianhui Liu, Siqi Guo, Yuming Lin, and Yong Li. Citygpt: Empowering urban spatial cognition of large language models. *arXiv preprint arXiv:2406.13948*, 2024.
- Leandro Augusto Frata Fernandes and Manuel Menezes de Oliveira. Geometric algebra: a powerful tool for solving geometric problems in visual computing. In 2009 Tutorials of the XXII Brazilian Symposium on Computer Graphics and Image Processing, pp. 17–30. IEEE, 2009.
- 627 Christian Freksa. Qualitative spatial reasoning. In *Cognitive and linguistic aspects of geographic* 628 *space*, pp. 361–372. Springer, 1991.
- Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making LLaMA SEE and draw with SEED tokenizer. In *The Twelfth International Conference on Learn- ing Representations*, 2024. URL https://openreview.net/forum?id=0Nui91LBQS.
- GeminiTeam, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
 capable multimodal models. *arXiv: 2312.11805*, 2023.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *CVPR*, pp. 6325–6334, 2017.
 - Joana Hois and Oliver Kutz. Towards linguistically-grounded spatial logics. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2011.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022.
- Michael Igorevich Ivanitskiy, Rusheb Shah, Alex F Spies, Tilman Räuker, Dan Valentine, Can Rager, Lucia Quirke, Chris Mathwin, Guillaume Corlouer, Cecilia Diniz Behn, et al. A configurable library for generating and manipulating maze datasets. *arXiv preprint arXiv:2309.10498*, 2023.

648 649 650	Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In <i>ICML</i> , volume 139, pp. 4904–4916, 2021.
651 652 653 654	Xinke Jiang, Ruizhe Zhang, Yongxin Xu, Rihong Qiu, Yue Fang, Zhiyuan Wang, Jinyi Tang, Hongxin Ding, Xu Chu, Junfeng Zhao, et al. Think and retrieval: A hypothesis knowledge graph enhanced medical large language models. <i>arXiv preprint arXiv:2312.15883</i> , 2023.
655 656	Christoph Kamann and Carsten Rother. Benchmarking the robustness of semantic segmentation models. In <i>CVPR</i> , pp. 8828–8838, 2020.
658 659	Amita Kamath, Jack Hessel, and Kai-Wei Chang. What's" up" with vision-language models? investigating their struggle with spatial reasoning. <i>arXiv preprint arXiv:2310.19785</i> , 2023.
660 661 662	Parisa Kordjamshidi, Martijn Van Otterlo, and Marie-Francine Moens. Spatial role labeling: To- wards extraction of spatial relations from natural language. <i>ACM Transactions on Speech and</i> <i>Language Processing (TSLP)</i> , 8(3):1–36, 2011.
664 665	Patrick Lester. A* pathfinding for beginners. <i>online]. GameDev WebSite. http://www. gamedev. net/reference/articles/article2003. asp (Acesso em 08/02/2009)</i> , 2005.
666 667 668	Chengzu Li, Caiqi Zhang, Han Zhou, Nigel Collier, Anna Korhonen, and Ivan Vulić. Topviewrs: Vision-language models as top-view spatial reasoners. <i>arXiv preprint arXiv:2406.02537</i> , 2024a.
669 670 671	Fangjun Li, David C Hogg, and Anthony G Cohn. Advancing spatial reasoning in large language models: An in-depth evaluation and enhancement using the stepgame benchmark. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 18500–18507, 2024b.
672 673 674 675	Hao Li, Xue Yang, Zhaokai Wang, Xizhou Zhu, Jie Zhou, Yu Qiao, Xiaogang Wang, Hongsheng Li, Lewei Lu, and Jifeng Dai. Auto mc-reward: Automated dense reward design with large language models for minecraft. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and</i> <i>Pattern Recognition</i> , pp. 16426–16435, 2024c.
676 677 678 679	Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. In <i>ICLR</i> , volume 162, pp. 12888–12900, 2022.
680 681 682	Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language- image pre-training with frozen image encoders and large language models. In <i>ICML</i> , volume 202, pp. 19730–19742, 2023a.
683 684 685 686 687	Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16</i> , pp. 121–137. Springer, 2020.
688 689 690	Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 23390–23400, 2023b.
691 692 693	Yuan-Hong Liao, Rafid Mahmood, Sanja Fidler, and David Acuna. Reasoning paths with reference objects elicit quantitative spatial reasoning in large vision-language models. <i>arXiv preprint arXiv:2409.09788</i> , 2024.
695 696	Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. <i>Transactions of the Association for Computational Linguistics</i> , 11:635–651, 2023a.
697 698 699	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In <i>NeurIPS</i> , 2023b.
700 701	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. URL https: //llava-vl.github.io/blog/2024-01-30-llava-next/.

702 703 704	Roshanak Mirzaee, Hossein Rajaby Faghihi, Qiang Ning, and Parisa Kordjmashidi. Spartqa:: A textual question answering benchmark for spatial reasoning. <i>arXiv preprint arXiv:2104.05832</i> , 2021.
705 706 707	OpenAI. Gpt-4v(ision) system card . https://cdn.openai.com/papers/GPTV_System_ Card.pdf, 2023.
708 709 710 711	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar- wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In <i>ICML</i> , volume 139, pp. 8748–8763, 2021a.
712 713 714 715 716	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021b.
710 717 718	Navid Rajabi and Jana Kosecka. Towards grounded visual spatial reasoning in multi-modal vision language models. <i>arXiv preprint arXiv:2308.09778</i> , 2023.
719 720	Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text- conditional image generation with clip latents. <i>arXiv preprint arXiv:2204.06125</i> , 1(2):3, 2022.
721 722 723 724	Kanchana Ranasinghe, Satya Narayan Shukla, Omid Poursaeed, Michael S Ryoo, and Tsung-Yu Lin. Learning to localize objects improves spatial reasoning in visual-llms. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 12977–12987, 2024.
725 726 727	Julia Rozanova, Deborah Ferreira, Krishna Dubba, Weiwei Cheng, Dell Zhang, and Andre Freitas. Grounding natural language instructions: Can large language models capture spatial information? <i>arXiv preprint arXiv:2109.08634</i> , 2021.
728 729 730 731	Dhruv Shah, Błażej Osiński, Sergey Levine, et al. Lm-nav: Robotic navigation with large pre- trained models of language, vision, and action. In <i>Conference on robot learning</i> , pp. 492–504. PMLR, 2023.
732 733 734	Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. Stepgame: A new benchmark for robust multi- hop spatial reasoning in texts. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 36, pp. 11321–11329, 2022.
735 736 737	Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension. In <i>ECCV</i> , pp. 742–758, 2020.
738 739	Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards VQA models that can read. In <i>CVPR</i> , 2019.
740 741 742	Yihong Tang, Zhaokai Wang, Ao Qu, Yihao Yan, Kebing Hou, Dingyi Zhuang, Xiaotong Guo, Jinhua Zhao, Zhan Zhao, and Wei Ma. Synergizing spatial optimization with large language models for open-domain urban itinerary planning. <i>arXiv preprint arXiv:2402.07204</i> , 2024.
743 744 745	Marc van Zee. Measuring compositional generalization, 2020. URL https://research.google/blog/measuring-compositional-generalization/.
746 747 748	Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. <i>arXiv preprint arXiv:2406.14852</i> , 2024.
749 750 751 752 753	Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In <i>International conference on machine learning</i> , pp. 23318–23340. PMLR, 2022a.
754	Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal,

 Weinid Walg, Hargoo Bao, Er Dong, Johan Björck, Zinnang Feng, Grang Eld, Kitti Aggatwal,
 Owais Khan Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv: 2208.10442*, 2022b.

100	Zhaokai Wang, Renda Bao, Qi Wu, and Si Liu. Confidence-aware non-repetitive multimodal trans-
757	formers for textcaps. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35,
758	pp. 2835–2843, 2021.
759	

- Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei.
 Visualization-of-thought elicits spatial reasoning in large language models. arXiv preprint arXiv:2404.03622, 2024.
- Jinxuan Xu, Shiyu Jin, Yutian Lei, Yuqian Zhang, and Liangjun Zhang. Reasoning tuning grasp:
 Adapting multi-modal large language models for robotic grasping. In 2nd Workshop on Language
 and Robot Learning: Language as Grounding, 2023.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of Imms: Preliminary explorations with gpt-4v (ision). *arXiv: 2309.17421*, 9, 2023a.
- Zhun Yang, Adam Ishay, and Joohyung Lee. Coupling large language models with logic programming for robust and general reasoning from text. *arXiv preprint arXiv:2307.07696*, 2023b.
- Taiping Zeng and Bailu Si. Cognitive mapping based on conjunctive representations of space and movement. *Frontiers in neurorobotics*, 11:61, 2017.
- Yichi Zhang, Yao Huang, Yitong Sun, Chang Liu, Zhe Zhao, Zhengwei Fang, Yifan Wang, Huanran Chen, Xiao Yang, Xingxing Wei, et al. Benchmarking trustworthiness of multimodal large
 language models: A comprehensive study. *arXiv preprint arXiv:2406.07057*, 2024.
- Yuze Zhao, Jintao Huang, Jinghan Hu, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. Swift: A scalable lightweight infrastructure for fine-tuning. *arXiv preprint arXiv:2408.05517*, 2024.

810 A APPENDIX

812 A.1 IMPLEMENTATION DETAILS

In addition to the experimental settings outlined in Section 4.1, we provide the following categorized
 implementation details for this work.

For model specifications, the GPT-40 model used in our experiments and demonstrations is based on the gpt-40-2024-05-13 version, while the Gemini model is Gemini 1.5 Flash. For TSP data generation, we used an open-source Python TSP solver¹ to obtain the ground truth visiting order of the given object coordinates.

For VLM evaluations, we focused on four directional categories (top left, top right, bottom left, and bottom right) to make it easier for VLMs to distinguish between directions. To discretize object locations for localization learning in VLM, the 2D space is proportionally divided using 40% and 60% thresholds along both the x and y axes, creating nine distinct regions (center, top, bottom, left, right, top-left, top-right, bottom-left, bottom-right). Detailed data statistics and distribution visualizations are provided in Section A.2.

To extract and format the VLMs' responses, we used the LLaMA 3.1 language model (Dubey et al., 827 2024), which converts the results into the required format for metric calculations. The specific 828 prompts used for each task are detailed in Section A.4. The evaluation for basic spatial relationship 829 understanding is intuitive, as it follows a multiple-choice question format. For the SPP evaluation, 830 we check two criteria: (1) whether the solved path is valid on the grid, and (2) whether the length of 831 the solved path is indeed the shortest between the given start and end objects. For the TSP evaluation, 832 a path is considered "correct" only if it exactly matches the solution from the TSP solver mentioned 833 above. To reduce the difficulty for VLMs in solving TSP, we explicitly specify the starting object in 834 our implementation.



A.2 DATA STATISTICS

835 836

858

859 860

861

862 863

Figure 7: Data statistics of basic spatial relationships (from left to right: distance, direction, and localization statistics).

To complement Section 4.1, this section provides detailed statistics of the data from Sparkle training set and evaluation. We begin by discussing data related to basic spatial relationships (i.e., distance,

¹https://github.com/fillipe-gsm/python-tsp

direction, localization), covering both Sparkle training set and the spatial relationship understanding task in the evaluation set.

Figure 7 illustrates various statistics. In the left column, we see the distribution of questions and instructions related to the *distance* between objects, which includes comparative expressions (e.g., shortest, shorter, longer, longest) and numerical distance estimations considered only in Sparkle training set. The training set shows a fairly even distribution of comparison queries, while in the test set, queries involving the "shortest" and "longest" distances occur more frequently than those involving "shorter" and "longer".

The middle column of Figure 7 presents the data concerning *directional* relationships between objects. We divided the 2D space into direction sectors: four sectors for testing and eight for training. The directional relationships of "bottom-right", "bottom-left", "top-right", and "top-left" each make up about 19% of the training data, while "top", "bottom", "left", and "right" each account for roughly 6%. In the test set, the four main directional relationships are distributed evenly.

Lastly, the right column in Figure 7 shows the *localization* data. Objects are most frequently located in the corners of the space (i.e., top-left, top-right, bottom-left, and bottom-right) in both the training and test sets. The number of objects placed in "top", "bottom", "left", and "right" positions is about half that of those in the corners, while the fewest objects are placed in the center. This is due to the intentional narrowing of the center area as we explained in Section A.1, which reduces the likelihood of randomly generated objects being placed there. Since there is no clear distinction between regions like "left" and "top-left", this narrowed design encourages VLMs to accurately distinguish specific areas such as the "center", "top", "bottom", "left", and "right"



Figure 8: Data statistics of composite spatial reasoning tasks in the evaluation set.

Figure 8 presents data statistics for composite spatial reasoning tasks. The two left subfigures show the distribution of ground truth shortest path lengths in 4×4 and 5×5 grids, while the two right subfigures depict the distribution of total distances for the optimal path in the TSP with 4 and 5 objects.

918 A.3 ADDITIONAL EXPERIMENTS

A.3.1 QWEN-VL EXPERIMENTS

	Spatia	l Relatio	nships	Traveling Sales		esman Problem		COCO-Spatial		GQA-Spatial		What's
Model	Loc.	Dist.	Dir.	4 Ob	ojects	5 Ob	jects	1 Object	2 Objects	1 Object	2 Objects	Up
	Acc.	Acc.	Acc.	Acc.	Marg.	Acc.	Marg.	Acc.	Acc.	Acc.	Acc.	Acc.
Qwen-VL-7B	22.0	37.0	24.5	10.5	0.126	2.5	0.458	89.8	74.3	98.5	94.0	42.7
+ Sparkle-Instruct	56.0	54.5	61.0	16.5	0.069	9.0	0.367	96.3	86.8	98.5	96.2	48.1
Δ	+155%	+47%	+149%	+57%	-45%	+260%	-20%	+7%	+17%	-	+2%	+13%

Table 5: Results of Qwen-VL Enhanced with Sparkle-Instruct

To further validate the generalizability of Sparkle, we conducted experiments using Qwen-VL (Bai et al., 2023b). The results, presented in Table 5, show significant improvements after fine-tuning with Sparkle compared to the original InternVL2-8B model.

Specifically, there is an approximately 120% improvement in the basic spatial relationship understanding task, a 150% improvement in the accuracy of composite tasks, and an 8% improvement in general spatial tasks. However, we excluded the results of the SPP task from our analysis, as the original performance of Qwen-VL-7B was too poor to allow for insightful comparison. This underperformance is not primarily due to limitations in spatial reasoning, but rather issues with visual recognition capabilities, as discussed in Section 4.2.1.

A.3.2 ABLATION STUDY RESULTS ON MARGIN METRIC



Figure 9: Margin results of InternVL2-8B with varying training sample sizes.

We have included the ablation study results on the margin metric here. As shown in Figure 9, a similar trend can be observed, consistent with the analysis discussed in Section 4.3.

A.4 PROMPTS FOR EXTRACTING INFERENCE RESULTS FROM VLMS

In this section, we provide the designed prompts for a LM to extract results from VLMs' responses.

A.4.1 MULTI-CHOICE QUESTIONS

Prompt for Extracting Results from VLMs' Responses to Multiple-Choice Questions X is the letter. Provide no additional content. The result is: ```{result}```.

The above prompt is adopted for all evaluations that in a Multi-choice Questions format.

A.4.2 SHORTEST PATH PROBLEM

```
Prompt for Extracting Results from VLMs' Responses to Shortest Path Problems
Extract the sequence of node labels from the given input and return it as a Python
list.
**Return Format:**
- Do not include any additional text or explanations.
- Ensure that the response is a single list containing only the node text labels (N1,
N2, ...).
- If no valid action sequence is found, return 'None'.
**Example Output format:**
[node1 text label, node2 text label, ...]
Now, extract the result from the following input: ```{result}```. Strictly adhere to
```

A.4.3 TRAVELING SALESMAN PROBLEM

the return format.

```
Prompt for Extracting Results from VLMs' Responses to Traveling Salesman Problems
Extract the sequence of movements from the given input and return it as a Python list
of object names.
**Return Format:**
- Do not include any additional text or explanations.
- Ensure that the response is a single list containing only the object names.
**Expected Output Format:**
{output_format}
Now, extract the result from the following response: "``{result}```. Strictly adhere to
the output format.
```

1026 A.5 SAMPLE DATA DEMONSTRATION

¹⁰²⁸ In this part, we provide detailed data sample from our experiments.

A.5.1 DATA SAMPLE FROM SPARKLE TRAINING SET AND EVALUATION



Figure 10: A data sample from the Sparkle training set.

A.5.2 DATA SAMPLE FROM THE BASIC SPATIAL RELATIONSHIP UNDERSTANDING TASK

tionship Understanding
Distance: Q: Which distance is the shortest? Options: A. N1 to N4, B. N1 to N3, C. N4 to N3
<pre>Direction: Q: Determine the direction from N1 to N2. Options: A. top left, B. top right, C. down left, D. down right</pre>
<pre>Localization: Q: What is the location of the N4 object? Options: A. top left, B. top, C. top right, D. left, E. center, F. right, G. bottom left, H. bottom, I. bottom right</pre>
A: A (Distance), B (Direction), I (Localization)





1080 A.5.3 DATA SAMPLE FROM THE SHORTEST PATH PROBLEM

	CALLASINGIO OF CHAT WITH	,	
A.6.1	GPT-40		
	ChatGPT 4o 🗸		,†, Share
			_
			N4 N8 N12 N16
			N3 N7 N11 N15
			N2 N6 N10 N14
			N1 N9
		The image shows a	grid graph where each node is labeled (N1,
		N2, N16) and cor image, find the sho	nnected to neighboring nodes. Based on the prtest path from the start node (green) to the
		end node (red) with	hout loops or backtracking. Don't use any
		shortest path in a li	any code^^. Directly return the solved ist format. Example output format: ["Na",
		"Nb", "Nc"].	
	You're gi	ving feedback on a	new version of ChatGPT.
	Which respon	nse do you prefer? Respo	onses may take a moment to load.
	Response 1		S Response 2
	The green node is labeled "N5" (start no	ode), and the	The shortest path from the green node (N5) to the re
	red node is labeled "N13" (end node). N the shortest path from N5 to N13 on th	low let's find is grid graph:	node (N13) without loops or backtracking can be determined by following the grid structure and
	1. N5 (start) \rightarrow N6		considering the direct connections between the
	2. N6 → N10		nodes.
	3. N10 → N9		Here is the path:
	4. N9 → N13 (end)		• Start at N5 (green)
	So, the shortest path from N5 to N13 is	:	Move to N6
			Move to N10
	[10 , 10 , 10 , 10 , 10].		End at N13 (red)
	I prefer this response		Thus, the shortest path is:
			["N5", "N6", "N10", "N13"].
			I prefer this response
	Figure 15: Screenshot supp	orting Figure	e 1: Chat interactions with GPT-40



1242 A.6.2 INTERNVL2-PRO

InternVL2

Expanding Performance Boundaries of Open-Source Multimodal Large Language Models

The image shows a grid graph where each node is labeled (N1, N2, ... N16) and connected to neighboring nodes. Based on the image, find the shortest path from the start node (green) to the end node (red) without loops or backtracking. Don't use any tools, **don't write any code**. Directly return the solved shortest path in a list format. Example output format: ["Na", "Nb", "Nc"].



(In this conversation, 1 image was uploaded, 1 image in total)

To determine the shortest path from the green start node (N5) to the red end node (N13), we need to traverse through the graph using only adjacent nodes.

Here's one possible path:

- Start at N5
- Move right to N6
- Move down to N10
- Move right to N11
- Move down to N15
- Move right to N14
- Move up to N13

Therefore, the shortest path can be represented as:

```
["N5", "N6", "N10", "N11", "N15", "N14", "N13"]
```

This path ensures that no node is visited more than once and reaches the destination efficiently.

Clear Chat History Regenerate

Figure 17: Screenshot supporting Figure 1: Chat interactions with InternVL2-Pro.