TOWARDS GEOMETRY PROBLEMS SOLVING EMPLOY-ING GPT-4 VISION WITH FEW-SHOT PROMPTING: AN EMPIRICAL STUDY OF WHAT MATTERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Few-shot prompting methods can significantly improve the ability of Large Language Models (LLMs) in mathematical reasoning, including geometry problem solving (GPS). GPT-4 Vision (GPT-4V), as a leading example of LLMs, also demonstrates significant improvements. This tremendous achievement is mainly attributed to prompting methods like "Chain-of-Thought" and "Programof-Thought," which leverage the in-context learning ability of the model combined with few-shot prompting to solve new problems. Despite the success of these prompting methods, it remains understood what the GPT-4V model learns from the demonstrations that leads to improved performance. In this paper, we evaluated the answering accuracy of GPT-4V with few-shot prompting on five geometric problem datasets and conducted a series of detailed analyses. Firstly, through ablation experiments with valid and invalid demonstrations, we found that the model's performance improvement is not due to the quality of the demonstration, but rather to the input format, output format, and logic and structure of the demonstration. Secondly, by analyzing the reasoning and computational requirements of GPS, and verifying experimental results, we found that GPS tasks emphasize reasoning ability more than computational power. Finally, our analysis of various prompt methods revealed that existing approaches are not effective at improving model performance concerning problem length and geometric shape. Therefore, specialized prompt methods could be designed to enhance the model's performance in these aspects, or fine-tuning the model by adding geometric problem data with longer lengths or mixed geometric shapes could optimize its performance. Overall, developing an LLM that fully adapts to GPS tasks represents a key research direction. The source code will be made available in a GitHub repository.

036

038

005 006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

030

031

032

033

034

1 INTRODUCTION

039 It is widely consensus that leveraging the reasoning and in-context learning capabilities of large 040 language models (LLMs) Zhang et al. (2023a); Kandpal et al. (2023); Shi et al. (2023); Ye et al. 041 (2023b), combined with few-shot prompting, can significantly improve their performance in math-042 ematical reasoning Yin et al. (2023); Gao & Das (2024); Firdaus et al. (2023); Liu et al. (2023); 043 Wu et al. (2023b). In mathematical reasoning research, geometry problem solving (GPS) Gao et al. 044 (2023a); Chang et al. (2022); Ning et al. (2023); Peng et al. (2023); Sun et al. (2024a) is crucial as it demands higher levels of specialized knowledge and comprehensive skills Lu et al. (2024), showcasing the potential of large language models more effectively Zhang et al. (2024). Therefore, 046 exploring methods to solve geometry problems using LLMs with few-shot prompting, along with 047 an in-depth analysis of key factors, can provide essential guidance and insights for future research 048 in the GPS field. 049

Currently, we know very little about what LLMs have learned from prompting methods Chen et al. (2023b); Wu et al. (2023a); Wang et al. (2023a); Gao et al. (2023c); Hu et al. (2024). The most direct way for LLMs to improve their ability to solve geometric problems is by providing a small number of examples, which prompts the model to answer new questions based on these examples. The current prompt methods are mainly divided into two categories: "Chain-of-Thought" Wei et al. (2022);

Zhou et al. (2023); Jiang et al. (2023); Wang et al. (2023b); Mondal et al. (2024) and "Program-of-055 Thought" Chen et al. (2023a); Gao et al. (2023b); Ye et al. (2023a); Imani et al. (2023). The "Chain-056 of-Thought" methods utilize the model's extensive pre-training parameters ('memory') Li et al. 057 (2023); Zhong et al. (2024) combined with provided demonstrations to guide the model in reasoning 058 the answer step by step Gao et al. (2023b). In contrast, the "Program-of-Thought" methods leverage the model's code generation ability to use external tools for performing complex computations Chen et al. (2023a). While these methods have shown some success in enhancing the model's problem-060 solving ability, the key factors driving this improvement remain unknown. To deeply analyze the 061 auxiliary effect of prompting methods on LLMs, we research several key issues: Firstly, are valid 062 demonstrations of "Chain-of-Thought" reasoning or "Program-of-Thought" computation important 063 for guiding GPT-4V in GPS tasks? Secondly, are the "Program-of-Thought" computation methods 064 superior to the "Chain-of-Thought" reasoning methods in GPS tasks? Finally, what other aspects 065 need to be considered when using prompting methods to solve geometric problems, beyond the 066 reasoning process and computation methods?

067 We have conducted a series of detailed experimental designs and analyses to address the aforemen-068 tioned key issues Wang et al. (2023a); Gao et al. (2023c). First, to evaluate the impact of including 069 only valid demonstrations on model performance, we conducted ablation experiments using eight different prompting methods. Meanwhile, each method included invalid demonstrations, allowing 071 us to assess the model's ability to learn between valid and invalid demonstrations. From the experi-072 mental results, we found that the improvement in the model's performance is not related to the value 073 of the demonstration, but rather to the input format, output format, and the logic and structure of the 074 demonstration. Second, to evaluate whether solving geometric problems is important for improving 075 reasoning or computational abilities, we analyzed the reasoning and computational requirements from the domain knowledge and the complexity of the geometric problems themselves. By com-076 paring the experimental results of two types of prompting methods, we found that solving most 077 geometric problems requires a stronger emphasis on reasoning ability, while only a few geomet-078 ric problems require complex computational power. Therefore, in GPS tasks, "Chain-of-Thought" 079 reasoning methods are superior to "Program-of-Thought" computation methods. Finally, we an-080 alyzed the effect of various prompting methods on the problem length and the geometric shapes 081 involved. The analysis of experimental results showed that these prompting methods did not significantly improve these aspects. This indicates that specialized prompting methods could be designed 083 to improve the model's performance in these aspects, or fine-tuning the model can be achieved by 084 adding problem data with longer lengths or mixed geometric shapes to optimize its performance. 085

Overall, our research and analysis of LLMs with few-shot prompting to solve geometric problems revealed that the model can learn the prompting framework, including the input format, output format, and answering ideas, from the demonstrations. However, the specific answer process still depends on the problem's content and relies on "memory" to generate an answer. Therefore, proposing a large language model that fully adapts to GPS tasks is a key research direction.

090 091

2 RELATED WORK

092 093 094

095

2.1 CHAIN-OF-THOUGHT METHODS

The Chain-of-Thought methods involve a series of reasoning prompts that divide a problem into mul-096 tiple intermediate steps, gradually solving each step to obtain the final answer ultimately. Among 097 these methods, chain-of-thought prompting (CoT) Wei et al. (2022), as shown in Figure 1 (a), 098 has achieved remarkable results in solving general reasoning tasks. However, its performance diminishes when tackling problems with higher reasoning difficulty. To address this issue, the least-100 to-most prompting (LtM) Zhou et al. (2023) method is proposed. LtM decomposes a complex 101 problem into a series of easier subproblems and solves them sequentially, using the answers from 102 previous subproblems to assist in solving subsequent ones. To clarify the solution process, the 103 plan-and-solve prompting (PS) method Wang et al. (2023b) is proposed. PS solves new problems 104 through a series of simple trigger sentences (such as 'give', 'plan', 'calculation', and 'answer'). 105 The standard CoT method follows a roughly linear reasoning approach but often forgets previous intermediate results during the process. The residual connection prompting (RP) method Jiang 106 et al. (2023) mitigates this issue by reintroducing the results of previous steps as prerequisites for 107 subsequent steps, thereby reducing the model's tendency to forget.



Figure 1: Comparison of the implementation process of two types of prompting methods. (a) Chainof-Thought methods rely entirely on the GPT-4V' "memory" for step-by-step reasoning and calculation; (b) Program-of-Thought methods generate programs and use external tools for more precise computation.

2.2 PROGRAM-OF-THOUGHT METHODS

The Program-of-Thought methods involve a series of computational prompts that utilize the pow-erful code-generation capabilities of LLMs to solve complex computational problems. Although LLMs excel at gradually decomposing problems using Chain-of-Thought methods, they often en-counter logical and arithmetic errors when solving individual subproblems, even if the overall decomposition is correct. To overcome these issues, two representative methods have been proposed: program-of-thought (PoT) Chen et al. (2023a) and program-aided language (PAL) Gao et al. (2023b), as shown in Figure 1 (b). Both methods separate reasoning from computation by placing each computational step into an external code executor (such as a Python interpreter) for execu-tion. The key difference is that PoT describes each step entirely in the programming language, while PAL integrates both natural language and programming language. To improve the effective-ness in solving search constraint problems, the satisfiability-aided language (SATLM) Ye et al. (2023a) is proposed. This method converts natural language reasoning problems into satisfiability (SAT) problems and then obtains answers by using an SAT solver. To enhance the credibility of LLMs-generated answers, MathPromoter (MP) Imani et al. (2023) uses hint techniques to gener-ate multiple algebraic expressions and Python functions, solving the same mathematical problem with different approaches. The consensus among these solutions serves as the final answer, thereby increasing the confidence in the output.

BACKGROUND & STUDY FORMULATION

There is a wide consensus that Language Models are Few-Shot LearnersBrown et al. (2020); Wu et al. (2023a), and GPT-4V is no exception Sun et al. (2024b); Jin et al. (2024), as reflected in nu-merous studies. In this section, we first defined two types of prompting methods, Chain-of-Thought Reasoning and Program-of-Thought Computation, and then detailed the input format, output format, and the logic and structure of the demonstration.

3.1 CHAIN-OF-THOUGHT REASONING

In the realm of ultra-large-scale unsupervised deep learning, large models are often perceived as black boxes Rai et al. (2023), with their reasoning and decision-making processes being difficult to explain Wang et al. (2023c). This lack of transparency poses a challenge to the credibility of the model results. However, the introduction of the few-shot chain-of-thought method offers a potential

solution Wei et al. (2022); Zhou et al. (2023); Jiang et al. (2023); Wang et al. (2023b); Mondal et al. (2024). This method breaks down logical reasoning problems into multiple steps, generating results with a clear logical context. This approach improves the interpretability of the model, allowing humans to understand the derivation process of the answer. Thus, the chain-of-thought reasoning method marks a crucial milestone in enhancing the success of LLMs.

Specifically, few-shot chain-of-thought reasoning involves LLMs solving a new problem p by following few-shot demonstrations $\langle p_k \cdot R_k \cdot a_k \rangle$ and generating an answer a with a reasoning process R. The standard implementation process is as follows:

$$GPT\left\{\left[S, \left(\langle p_1 \cdot R_1 \cdot a_1 \rangle, \cdots, \langle p_k \cdot R_k \cdot a_k \rangle\right)\right] \oplus p\right\} \to \langle R, a \rangle$$

$$\tag{1}$$

where S is a system prompt used to guide the model in generating specific information or completing specific tasks when answering problems. In (a) Chain-of-Thought of Figure 1, we used the system prompt *Solve the question through step-by-step reasoning following the given examples*. Additionally, k represents the number of demonstrations, typically an integer not less than 0. When k = 0, it indicates the zero-shot prompting method; when k = 1, it indicates the one-shot prompting method; and when $k \ge 2$, it indicates the few-shot prompting method.

$$RE\left(\langle R, a \rangle\right) \to a \tag{2}$$

¹⁸³ Due to the concatenation of reasoning and the answer in the results generated by the LLM, we generally use Regular Expressions (*RE*) to extract the pure answer. For example, in (a) Chain-of-Thought of Figure 1, the output *R* is: "Since line AC=6in, and this is a regular polygon, therefore, the perimeter of the regular polygon is 6+6+6=18in. The answer is 18." Here, the answer a is 18, and the regular expression we used is: r "The answer is $(\backslash d+)$ ".

189 3.2 PROGRAM-OF-THOUGHT COMPUTATION190

171

172 173

180 181 182

188

201 202 203

209 210 211

Even with chain-of-thought reasoning, LLMs do not truly understand mathematical logic or the
fundamental concepts of addition, subtraction, multiplication, and division Zhong et al. (2024);
Zhou et al. (2024). Instead, they rely on prior knowledge to mimic problem-solving processes,
much like "drawing a dipper with a gourd as a model" (a Chinese idiom meaning to imitate without
understanding). Therefore, for tasks requiring precise arithmetic, professional computing tools are
still necessary, leading to the development of program-of-thought computing Chen et al. (2023a);
Gao et al. (2023b); Ye et al. (2023a); Imani et al. (2023).

Specifically, few-shot program-of-thought computation is a method where LLMs solve a new problem p by following few-shot demonstrations $< pk \cdot Ck >$ and generating an answer a using the program C they generate. The standard implementation process is as follows:

$$GPT \{ [S, (\langle p_1 \cdot C_1 \rangle, \dots \langle p_k \cdot C_k \rangle)] \oplus p \} \to C$$
(3)

where S is a system prompt designed to guide the model in generating specific information or completing particular tasks when solving problems. In (b) Program-of-Thought of Figure 1, the system prompt *Generate a code block following the given examples to solve the question* is used. Here, k represents the number of demonstrations provided.

$$Interpreter\left(C\right) \to a \tag{4}$$

The core idea of the program-of-thought method for achieving precise arithmetic computation is to
input the generated program or code block C into a program interpreter, using professional tools
for high-precision computation. For example, in (b) Program-of-Thought of Figure 1, the output C
is a *solution* function code block written in Python. We used a Python interpreter to execute this function, obtaining the final answer 18.

Table 1: Details of datasets being tested. The "total" represents the problem number of questions in
an original dataset, and the "sample" represents the number of problems sampled from a dataset in
a test.

Dataset	Total	Sample	Average problem words	Average knowledge
GEOS	186	62	24.7	1.3
Geometry3K	3002	1000	12.2	1.6
GeoQA	4998	1666	52.5	2.1
GeoQA+	7528	2510	54.5	1.8
PGPS9K	9022	2800	17.8	1.7

227 3.3 RESEARCH QUESTIONS

224 225 226

233

234

235

237

238

239

240

241 242

243 244

245

264

265

In the prompting methods discussed above, valid reasoning or program examples are provided as demonstrations to illustrate how GPT-4V derives the generated answer to a new problemWu et al. (2023b); Lu et al. (2024). Despite the impressive performance of various prompting methods in mathematical reasoning tasks, we are interested in exploring the following questions:

- Q1: Do valid demonstrations of chain-of-thought reasoning or program-of-thought computation matter for guiding GPT-4V in performing GPS tasks? If not, what does GPT-4V learn to obtain the answer?
- Q2: Is the program-of-thoughts computation superior to chain-of-thought reasoning prompting method in GPS task? If not, what characteristics in geometry problems would cause this phenomenon to occur?
 - Q3: What other aspects need to be considered when using prompting methods to solve geometry problems, besides the reasoning process and computation methods?
- 4 EXPERIMENTAL SETUP
- 4.1 DATASETS & IN-CONTEXT EXEMPLARS

246 The goal of our experiment is to analyze which factors are important for using prompting methods 247 with few-shot demonstrations to assist LLMs in solving geometric problems. Therefore, the dataset 248 used in our experiment includes five publicly available geometric problem datasets. The detailed 249 introduction is as follows: (1) GEOS Seo et al. (2015): the dataset contains simple middle school 250 geometry problems with geometric shapes. (2) Geometry3K Lu et al. (2021): the dataset contains 251 numerous geometry problems where semantic information is scarce and most values need to be 252 obtained from images. (3) GeoQA Chen et al. (2021): the dataset contains rich semantic information for middle and high school geometry problems. (4)GeoOA+ Cao & Xiao (2022): the dataset is based 253 on GeoQA, which adds more diverse types of geometry problems and forms an enhanced benchmark 254 dataset. (5) **PGPS9K** Zhang et al. (2023b): the dataset has both fine-grained graph annotations and 255 interpretable solution programs, and a small portion of the dataset comes from Geometry3K. Due 256 to budget considerations, we sample a certain number of questions from the five geometric problem 257 datasets being tested, and the number of samples is shown in Table 1. 258

Since our testing task only involved solving geometric problems, the five datasets share the same
problem-prompting template. The only difference is that the demonstrations of the solving process for the same problem are designed based on different prompt methods, including CoTWei
et al. (2022), LtMZhou et al. (2023), PSWang et al. (2023b), RPJiang et al. (2023), PoTChen et al.
(2023a), PALGao et al. (2023b), SATLMYe et al. (2023a), and MathprompterImani et al. (2023).

4.2 EVALUATION

We employed the GPT-4V API (gpt-4-turbo)¹, the mature GPT-4 Turbo model with vision capabilities, for our experiment. To evaluate the performance of GPT-4V with few-shot prompts, we used two strategies: (1) Average answering accuracy: we randomly sampled problems and tested the

¹https://platform.openai.com/docs/models/gpt-4-turbo-and-gpt-4

model for M rounds, taking the average accuracy as the performance metric. (2) Self-consistency answering accuracy: we evaluated performance with self-consistency decoding Wang et al. (2023d); Weng et al. (2023) by uniformly sampling and aggregating the most consistent answer from the Ncandidate answers. According to the usual research settings Gao et al. (2023b); Wang et al. (2023d); Weng et al. (2023); Imani et al. (2023), M and N are set to 40, and the temperature T is set to 0.

275 276 277

278

279

280

281

282

283

5 HOW MUCH DOES VALID DEMONSTRATION MATTER? (Q1)

Intuitively, one of the most critical aspects of a demonstration is its logical validity and sound reasoning. If we provide a demonstration with invalid reasoning steps and computational processes, we would expect GPT-4V to fail to answer properly and potentially experience little to no improvement, or even a decrease in accuracy, compared to standard answering accuracy. This is because we are guiding GPT-4V to answer in the wrong way, which could be detrimental rather than beneficial. To test this intuition, we propose a series of studies where we construct invalid demonstrations for various prompting methods and measure their influence on model behavior.

284 285 286

5.1 CONSTRUCTING INVALID DEMONSTRATIONS

We have set valid and invalid prompting demonstrations for eight prompting methods Wang et al.
(2023a), as shown in Appendix A. Specifically, valid demonstrations are the standard solving processes that can obtain a standard answer for the demonstration problem; invalid demonstrations are not relevant to the demonstration problem (including numerical and textual descriptions), and the standard answer cannot be obtained.

For example, the demonstration geometry problem in Appendix A is From the image, you can see 293 that the shorter base of the trapezoid is 6ft, the longer base is 16ft, and the height is 12ft. Find 294 the area of the trapezoid. For the representative prompting method in the Chain of Thought Rea-295 soning series, CoT prompting method, we have designed a valid solving demonstration for solving 296 the geometry problem based on the method: In this image, since the trapezoid has a longer base 297 (base1) with a length of 16 and a shorter base (base2) with a length of 6. The height, which is the 298 perpendicular distance between the two bases, is given as 12. To find the area of a trapezoid, you 299 can use the formula: Area = 1/2 * (base1 + base2) * height = 1/2 * (16 + 6) * 12=132. The area 300 of the trapezoid is 132. The answer is 132. Meanwhile, based on previous research, we have set up 301 a completely invalid demonstration, as shown in Appendix A.

302 303

304

5.2 RESULTS & ANALYSIS

305 **Results.** Figure 2 shows the answering accuracy of GPT-4V with one-shot demonstrations for solving geometric problems under different prompting settings (valid and invalid). From the comparison 306 of the valid-invalid experimental results, it can be seen that invalid prompting settings have both 307 an increase and a decrease compared to valid ones. The most significant increase is LtM with the 308 invalid prompting setting, in the Average answering accuracy of Geometry3K datasets, which in-309 creased by 7.2% compared to that with the valid. The most significant decrease is PoT with the 310 invalid prompting setting, in the Average answering accuracy of GeoQA+ datasets, which decreased 311 by 6.3% compared to that with the valid. More detailed experimental data is presented in Appendix 312 C. Additionally, we present GPT-4V's answer accuracy with one-shot demonstrations in Appendix 313 B and with three-shot demonstrations in Appendix D.

314 Analysis. On the one hand, through a comparative analysis of valid and invalid demonstrations, 315 we found that invalid prompting settings have both an increase and a decrease compared to valid 316 ones. This indicates that there is no correlation between the prompting validity and the answering 317 accuracy. In other words, valid demonstrations do not matter for GPT-4V with few-shot prompting 318 in GPS tasks. Moreover, this increase or decrease is not particularly significant, indicating that the 319 influence of valid and invalid demonstrations in guiding GPT-4V are consistent, and proving that 320 GPT-4V has learned the same content from both valid and invalid prompting settings, including 321 input format, output format, and logical thinking. On the other hand, compared to the standard GPT-4V, we found that various prompting methods have obvious improvements in GPS tasks. Moreover, 322 we found that the self-consistency strategy has a higher answering accuracy and a smaller increase 323 or decrease compared to the average evaluation strategy. This situation arises because the self-



Figure 2: Comparison of the answering accuracy of GPT-4V employing prompt methods for geometry problems between valid and invalid prompting demonstrations.

consistency strategy is also a prompting method, which guides GPT-4V to find consensus from multiple candidates' answers. This indicates that the self-consistency prompting method is more stable and reliable in the GPS task.

5.3 SUMMARY

352 Based on the experimental results and further analysis in this section, we have summarized two 353 points: in GPS tasks, firstly, the valid demonstrations do not matter for GPT-4V with few-shot prompting, but GPT-4V can learn input format, output format, and logical thinking from the demon-354 strations. Secondly, although the prompting method significantly improves the performance of GPT-355 4V, the effects achieved by various prompt methods are different. Among them, the self-consistency 356 strategy overcomes the instability of the model itself to obtain more accurate and reliable answers. 357 However, in GPS tasks, we need to consider more how to improve the reasoning and computing ca-358 pabilities of GPT-4V. Therefore, we need to further analyze a question: Is the "Program-of-Thought 359 Computation" superior to the "Chain-of-Thought Reasoning" prompting methods in GPS task? If 360 not, what characteristics in geometric problems would cause this phenomenon to occur? 361

362

342

343

344 345 346

347

348

349 350

351

363

WHICH A SERIES PROMPTING METHOD MATTERS MORE? (Q2) 6

364 We tested two types of prompting methods - "Chain-of-Thought Reasoning" and "Program-of-365 Thought Computation". Intuitively, in GPS tasks, "Chain of Thought Reasoning", which relies 366 solely on memory to solve problems, seems to have no advantage. Instead, the "Program of Thought Computing", which uses tools to enhance performance, can obtain more accurate answers. To test 368 this intuition, we design a series study where we have provided basic evidence for evaluating two 369 types of prompting methods by analyzing the reasoning and computational requirements in GPS 370 tasks.

371 372

373

367

6.1 STATISTICAL REQUIREMENTS

374 **Reasoning.** The reasoning requirements for solving a geometric problem are positively correlated 375 with the domain knowledge involved in the problem, so, the more domain knowledge the problem involves, the more reasoning steps are required. For example, the domain knowledge involved in 376 problem p in Figure 1 is only one - "Isosceles (Equilateral) Triangle", so the reasoning requirements 377 for the problem are not significant. From the analysis of the average domain knowledge (Average

378 knowledge) in Table 1, the average domain knowledge of all four datasets exceeded 1.5, indicating 379 the existence of geometric problems involving many domain knowledge, further proving that GPS 380 has a clear reasoning requirement. In Appendix E, we further refined the distribution of domain 381 knowledge involved in the problem. We analyzed and found that problems involving more than 2 382 domain knowledge accounts for a considerable proportion, with 71.4% in the GeoQA dataset and 40.9% in the GeoQA+ dataset. Specifically, in the GeoQA dataset, 9.4% of the questions involve 383 domain knowledge exceeding 4. This indicates that GPS tasks require complex reasoning to be 384 completed. 385

386 **Computation.** The computational requirements for solving geometric problems are related to multi-387 digit arithmetic. Multi-digit arithmetic refers to arithmetic operations involving numbers. For ex-388 ample, the numerical values involved in solving the problem p_5 in Appendix F include 222, and 38707.567, where the largest number 38707.567 is a 5-digit number and the smallest number 222 is 389 a 3-digit number. Therefore, we take 5 (the largest) as the multi-digit. Additionally, we consider the 390 arithmetic computation of decimals numerical values as 0-digit. For example, the numerical values 391 involved in solving the problem p_0 in Appendix F include $\frac{4}{7}$, $\frac{5}{7}$, and 0.429, which are regarded as 392 decimals. Therefore, we take 0 (the decimals) as the multi-digit. In Figure 3, we statistically an-393 alyzed the distribution of problems with different multi-digit arithmetic in five datasets. We found 394 that the biggest computational requirement in these datasets is also in the five-digit arithmetic, and 395 the vast majority (over 98%) of problems are lower than 3-digit arithmetic. Therefore, the GPS task 396 requires a small amount of computation. 397

6.2 RESULTS & ANALYSIS

398

399

400 **Results.** Figure 2 shows the answering ac-401 curacy of GPT-4V with two-shot prompting methods by two evaluation strategies (aver-402 age and self-consistency answering accuracy). 403 Firstly, two different background colors repre-404 sent different prompting methods: the white 405 background in the table represents the "Chain-406 of-Thought Reasoning" series methods (CoT, 407 LtM, PS, and RP), the gray background repre-408 sents the "Program-of-Thought Computation" 409 series methods (PoT, PAL, SATLM, and Math-410 prompter). By comparing the accuracy of these 411 two series methods, it can be seen that the for-412 mer has a significant advantage in GPS tasks. For example, in the evaluation of the average 413 answering accuracy, the RP method with in-414 valid reasoning (44.9%) on the GeoQA dataset 415 improved the accuracy by 22.3% compared 416 to the PAL method with invalid computation



Figure 3: Statistical analysis of computational requirements for different datasets.

417 (22.6%). Meanwhile, in the evaluation of the self-consistency answering accuracy, the RP method 418 with valid reasoning (53.1%) on the GEOS dataset improved the accuracy by 22.5% compared to 419 the PAL method with invalid computation (30.6%). This indicates that the method of enhancing rea-420 soning ability is more effective for GPS tasks than computation. Furthermore, we also found that the 421 SATLM and Mathprompt prompting methods (belonging to the "Program of Thought Calculation" 422 series methods) exceeded some of the "Chain-of-Thought Reasoning" series methods. For exam-423 ple, in the self-consistency accuracy evaluation of the PGPS9K dataset, the Mathprompt prompting method with invalid computation (40.4%) outperforms all "Chain-of-Thought Reasoning" methods. 424 In addition, by comparing the answering accuracy on two similar datasets (GeoQA and GeoQA+), 425 we found that the accuracy on the GeoQA dataset was lower than that on the GeoQA+ dataset, 426 whether it was the standard GPT (25.4% and 26.5%) or human performance (61.2% and 66.4%). 427 However, after using the prompting method, except for PoT and PAL methods, the accuracy in the 428 GeoQA dataset is generally higher than that in the GeoQA+ dataset. 429

Analysis. On the one hand, according to numerous research analyses, the "Program-of-Thought Computation" series methods are better at handling the problem with large numbers than the "Chainof-Thought Reasoning" methods, such as $134672 \times 98564=?$. But does the GPS task require this 432 ability? From the analysis in Figure 3, we find that the GPS task requires a small amount of compu-433 tation. Therefore, there is a phenomenon that the answering accuracy of the "Chain-of-Thought Rea-434 soning" methods is higher than the "Program-of-Thought Computation" series methods. However, it 435 is an exception in SATLM and Mathprompt promotion methods, as these two methods are different 436 from PoT and PAL methods (belonging to the "Program of Thought Calculation" series methods) that rely entirely on the programs for reasoning and computation. On the contrary, they separate 437 the reasoning and computation processes, and the two parts complement each other. Therefore, 438 SATLM and Mathprompt promotion methods surpass some of the "Chain-of-Thought Reasoning" 439 methods. On the other hand, as shown in Appendix E, most problems in the GeoQA dataset require 440 complex reasoning processes, therefore there is a high demand for reasoning ability. Moreover, in 441 the comparison between two similar datasets (GeoQA and GeoQA+) in experimental results, we 442 found that compared to the GeoQA+ dataset, most prompting methods showed a more significant 443 improvement in the GeoQA dataset. This indicates that these methods greatly cater to the complex 444 reasoning requirement of the GeoQA dataset. 445

6.3 SUMMARY

448 Based on the mutual verification between GPS task requirements and experimental results, we sum-449 marize as follows: Firstly, compared to computational requirements, the reasoning requirements 450 are higher. Therefore, the "Chain of Thought Reasoning" series methods (CoT, LtM, PS, and RP 451) and some "Program of Thought Computing" methods (SATLM and Mathprompter) that guide model reasoning cater to complex reasoning requirements and have more significant improvement 452 effects. Secondly, separating computation from reasoning and using the reasoning process to guide 453 precise computation is an optimal prompting method, such as the SATLM and Mathprompter meth-454 ods. Besides reasoning and computational requirements, what other factors can affect the effective-455 ness of solving geometric problems? 456

457 458

459

465

467

446

447

7 WHAT OTHER ASPECTS ALSO MATTER? (Q3)

To analyze whether other aspects besides reasoning
and computing requirements would affect the ability
of GPT-4V with few-shot prompting to solve geometric problems, we mainly completed two evaluations: geometry problem length and geometry shape.

466 7.1 GEOMETRY PROBLEM LENGTH

The length of geometric problems represents the 468 number of word tokens in the problem text. For ex-469 ample, in Appendix F, p_0 : find x. This problem con-470 tains two words, "find" and "x", and its length is 2. 471 Intuitively, as the problem length increases, the more 472 semantic information the model needs to understand, 473 the more difficult it is to answer the problem, and 474 the lower the accuracy of the answer. Conversely, 475 the higher the accuracy of the answer. However, in 476 GPS tasks, a lot of information is contained in geo-477 metric shapes, and the information contained in the text is limited, resulting in shorter problem lengths, 478 such as p_3 in Appendix F. So for solving geometric 479



Figure 4: Answering accuracy of GPT-4V with different prompting under different problem lengths, in the GeoQA dataset.

problems, we need to distinguish it from general mathematical reasoning tasks, provide a specific
 relationship between the problem length and the answering accuracy, and analyze which range of
 problem length can obtain the optimal answering accuracy.

The answering accuracy of different problem lengths is shown in Figure 4. The experimental results
 were obtained using the average answering accuracy as an evaluation strategy employing GPT-4V
 with two-shot PAL and COT at different problem lengths. In addition, to highlight the effectiveness of the prompting method, we also used the standard GPT-4V (without any prompting) as the base-

line. From Figure 4, we can see that regardless of whether there is prompting or not, the accuracy trend in answering at different problem lengths is consistent. And the higher accuracy is concentrated between problem lengths of (40, 50). This indicates that the problem length is unrelated to the method with or without prompting, but only to the model's ability to understand semantic information.

491 492

493 494

7.2 GEOMETRY SHAPE

The fundamental reason why geometric problems differ from general mathematical reasoning problems is that many geometric problems contain abstract geometric shapes, which also pose a huge challenge in GPS tasks - cognitive geometric shapes. To analyze the cognitive ability of GPT-4V towards different geometric shapes, we evaluated the answering accuracy of GPT-4V with different prompting under different geometric shapes in the Geometry3K dataset. The experimental results are shown in Figure 5.

The experimental results were obtained using the av-501 erage answering accuracy as an evaluation strategy 502 employing GPT-4V with 2-shot PAL and COT at 503 different geometry shapes. In addition, to highlight 504 the effectiveness of the prompting method, we also 505 used the standard GPT-4V (without any prompting) 506 as the baseline. From Figure 5, we can see that GPT-507 4V has strong cognitive abilities for shapes such 508 as squares, rectangles, and parallelograms. Among 509 them, the accuracy of the CoT prompting method 510 for answering problems involving parallelograms reached 28%. This indicates that the current GPT-511 4V has good cognitive abilities for simple geometric 512 shapes, but there is still a lot of room for improve-513 ment. Furthermore, there is no correlation between 514 the use of prompting methods and the improvement 515 of answering accuracy. 516



Figure 5: Answering accuracy of GPT-4V with different prompting under different geometry shapes, in the Geometry3K dataset.

7.3 SUMMARY

520 Based on the experimental analysis of geometric problem length and geometric shape, we summa-521 rize as follows: Firstly, compared to the standard GPT-4V, the existing prompting methods do not 522 significantly improve the accuracy of answering for a certain problem length or geometric shape. 523 This also provides a starting idea for our future innovative prompting methods. For example, we 524 provide a targeted prompting method for particularly long problems; Alternatively, we can provide 525 visual cues for a certain geometric shape to make it easier for the model to recognize the shape and enhance its cognitive effect. Secondly, the most important thing is to enhance the model's ability to 526 solve geometric problems. Fine-tuning methods are recommended to improve the performance of 527 the model while ensuring that the visual features of geometric shapes and the semantic information 528 of longer problem texts can be fully understood. 529

530

517 518

519

531 532

533

8 CONCLUSION

In this paper, we aim to better understand what GPT-4V has learned from the few-shot demonstrations, we conducted a series of experiments and detailed analysis. We find that: (1) The model's performance improvement is not due to the quality of the demonstration, but rather to the input format, output format, and the logic and structure of the demonstration; (2) GPS tasks emphasize reasoning ability more than computational power; (3) Specialized prompting methods could be designed to enhance the model's performance. Overall, developing an LLM that fully adapts to GPS tasks represents a key research direction.

540 REFERENCES 541

558

559

560

571

577

592

- 542 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, and Arvind Neelakantan. Language models are few-shot learners. In Advances in Neural Infor-543 mation Processing Systems 33: Annual Conference on Neural Information Processing Systems 544 2020 (NeurIPS-2020), 2020.
- 546 Jie Cao and Jing Xiao. An augmented benchmark dataset for geometric question answering through 547 dual parallel text encoding. In Proceedings of the 29th International Conference on Computa-548 tional Linguistics (COLING-2022), pp. 1511–1520, 2022. 549
- Tyler A. Chang, Zhuowen Tu, and Benjamin K. Bergen. The geometry of multilingual language 550 model representations. In Proceedings of the 2022 Conference on Empirical Methods in Natural 551 Language Processing (EMNLP-2022), pp. 119–136, 2022. 552
- 553 Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, and Lingbo Liu. Geoqa: A geometric 554 question answering benchmark towards multimodal numerical reasoning. In Findings of the As-555 sociation for Computational Linguistics (ACL-2021), volume ACL/IJCNLP 2021, pp. 513–523, 556 2021.
 - Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. Transactions on Machine Learning Research (TMLR-2023), 2023a.
- 561 Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, and Dayiheng Liu. Hallucination 562 detection: Robustly discerning reliable answers in large language models. In *Proceedings of the* 563 32nd ACM International Conference on Information and Knowledge Management (CIKM-2023), pp. 245–255, 2023b. 565
- 566 Mauzama Firdaus, Gopendra Vikram Singh, Asif Ekbal, and Pushpak Bhattacharyya. Multi-step prompting for few-shot emotion-grounded conversations. In Proceedings of the 32nd ACM Inter-567 national Conference on Information and Knowledge Management (CIKM-2023), pp. 3886–3891, 568 2023. 569
- 570 Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, and Lanqing Hong. G-llava: Solving geometric problem with multi-modal large language model. CoRR, 572 abs/2312.11370, 2023a. 573
- 574 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, and Pengfei Liu. PAL: program-aided language 575 models. In International Conference on Machine Learning (ICML-2023), volume 202, pp. 10764-10799, 2023b. 576
- Shuzheng Gao, Xin-Cheng Wen, Cuiyun Gao, Wenxuan Wang, Hongyu Zhang, and Michael R. 578 Lyu. What makes good in-context demonstrations for code intelligence tasks with llms? In 579 38th IEEE/ACM International Conference on Automated Software Engineering (IEEE-2023), pp. 580 761-773, 2023c. 581
- 582 Xiang Gao and Kamalika Das. Customizing language model responses with contrastive in-context 583 learning. In Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-2024), pp. 18039-18046, 2024. 584
- 585 Wenbo Hu, Yifan Xu, Yi Li, Weiyue Li, Zeyuan Chen, and Zhuowen Tu. Bliva: A simple multi-586 modal llm for better handling of text-rich visual questions. In Proceedings of Thirty-Eighth AAAI 587 Conference on Artificial Intelligence (AAAI-2024), pp. 2256–2264, 2024. 588
- 589 Shima Imani, Liang Du, and Harsh Shrivastava. Mathprompter: Mathematical reasoning using large 590 language models. In Proceedings of The 61st Annual Meeting of the Association for Computa-591 tional Linguistics: Industry Track (ACL-2023), pp. 37-42, 2023.
- Song Jiang, Zahra Shakeri, Aaron Chan, Maziar Sanjabi, and Hamed Firooz. Resprompt: Residual connection prompting advances multi-step reasoning in large language models. CoRR, 2023.

613

632

633

634

635

639

640

641

594	Oiao Jin, Fangyuan Chen, Yiliang Zhou, Ziyang Xu, Justin M. Cheung, Robert Chen, Ronald M.
595	Summers, and Justin F. Rousseau. Hidden flaws behind expert-level accuracy of gpt-4 vision in
596	medicine. CoRR, abs/2401.08396, 2024.
597	

- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language
 models struggle to learn long-tail knowledge. In *International Conference on Machine Learning* (*ICML-2023*), volume 202, pp. 15696–15707, 2023.
- Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix X.
 Yu, and Sanjiv Kumar. Large language models with controllable working memory. In *Findings* of the Association for Computational Linguistics (ACL-2023), pp. 1774–1793, 2023.
- Jingping Liu, Tao Chen, Zujie Liang, Haiyun Jiang, and Yanghua Xiao. Hierarchical prompt tuning
 for few-shot multi-task learning. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM-2023)*, pp. 1556–1565, 2023.
- Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-2021)*, pp. 6774– 6786, 2021.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, KaiWei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of
 foundation models in visual contexts. In *International Conference on Learning Representations*(*ICLR-2024*), 2024.
- Debjyoti Mondal, Suraj Modi, Subhadarshi Panda, Rituraj Singh, and Godawari Sudhakar Rao.
 Kam-cot: Knowledge augmented multimodal chain-of-thoughts reasoning. In *Proceedings of Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-2024)*, pp. 18798–18806, 2024.
- Maizhen Ning, Qiu-Feng Wang, Kaizhu Huang, and Xiaowei Huang. A symbolic characters aware
 model for solving geometry problems. In *Proceedings of the 31st ACM International Conference on Multimedia (ACM-2023)*, pp. 7767–7775, 2023.
- Shuai Peng, Di Fu, Yijun Liang, Liangcai Gao, and Zhi Tang. Geodrl: A self-learning framework for geometry problem solving using reinforcement learning in deductive reasoning. In *Findings of the Association for Computational Linguistics: ACL 2023 (ACL-2023)*, pp. 13468–13480, 2023.
- Daking Rai, Yilun Zhou, Bailin Wang, and Ziyu Yao. Explaining large language model-based neural
 semantic parsers. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence* (AAAI-2023), pp. 16308–16309, 2023.
 - Min Joon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. Solving geometry problems: Combining text and diagram interpretation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP-2015)*, pp. 1466–1476, 2015.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, and Ed H. Chi. Large
 language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning (ICML-2023)*, volume 202, pp. 31210–31227, 2023.
 - Kai Sun, Yushi Bai, and Nianyi Lin. Advancing geometric problem solving: A comprehensive benchmark for multimodal model evaluation. *CoRR*, abs/2404.05091, 2024a.
- Qi Sun, Xiao Cui, Wengang Zhou, and Houqiang Li. Exploiting gpt-4 vision for zero-shot point cloud understanding. *CoRR*, abs/2401.07572, 2024b.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun.
 Towards understanding chain-of-thought prompting: An empirical study of what matters. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL-2023)*, pp. 2717–2739, 2023a.

- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, and Yunshi Lan. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL-2023)*, pp. 2609–2634, 2023b.
- Kinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language models are latent variable models: Explaining and finding good demonstrations for incontext learning. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems (NeurIPS-2023), pp. 16308–16309, 2023c.
- Kuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha
 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
 models. In *Proceedings of The Eleventh International Conference on Learning Representations* (*ICLR-2023*), 2023d.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, and Brian Ichter. Chain-of-thought
 prompting elicits reasoning in large language models. In Advances in Neural Information Process *ing Systems 35: Annual Conference on Neural Information Processing Systems 2022 (NeurIPS-2022)*, 2022.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and
 Jun Zhao. Large language models are better reasoners with self-verification. In *Findings of the Association for Computational Linguistics (EMNLP-2023)*, pp. 2550–2575, 2023.
- Wenhao Wu, Huanjin Yao, Mengxi Zhang, Yuxin Song, Wanli Ouyang, and Jingdong Wang.
 Gpt4vis: What can GPT-4 do for zero-shot visual recognition? *CoRR*, abs/2311.15732, 2023a.
- Yiran Wu, Feiran Jia, Shaokun Zhang, Hangyu Li, Erkang Zhu, Yue Wang, Yin Tat Lee, Richard
 Peng, Qingyun Wu, and Chi Wang. An empirical study on challenging math problem solving
 with gpt-4. *CoRR*, abs/2306.01337, 2023b.
- Ki Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett. Satlm: Satisfiability-aided language models
 using declarative prompting. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023 (NeurIPS-2023)*, 2023a.
- Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. Large language models are versatile decomposers: Decomposing evidence and questions for table-based reasoning. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR-2023)*, pp. 174–184, 2023b.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *CoRR*, abs/2306.13549, 2023.
- Biao Zhang, Barry Haddow, and Alexandra Birch. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning (ICML-2023)*, volume 202, pp. 41092–41110, 2023a.
- Duzhen Zhang, Yahan Yu, Chenxing Li, Jiahua Dong, Dan Su, Chenhui Chu, and Dong Yu. Mm Recent advances in multimodal large language models. *CoRR*, abs/2401.13601, 2024.
- Mingliang Zhang, Fei Yin, and Cheng-Lin Liu. A multi-modal neural geometric solver with textual
 clauses parsed from diagram. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence (IJCAI-2023)*, pp. 3374–3382, 2023b.
- Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large
 language models with long-term memory. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-2024)*, pp. 19724–19731, 2024.
- Denny Zhou, Nathanael Scharli, Le Hou, Jason Wei, Nathan Scales, and Xuezhi Wang. Least-to most prompting enables complex reasoning in large language models. In *The Eleventh Interna- tional Conference on Learning Representations (ICLR-2023)*, 2023.
- Yujia Zhou, Qiannan Zhu, Jiajie Jin, and Zhicheng Dou. Cognitive personalized search integrating large language models with an efficient memory mechanism. In *Proceedings of the ACM on Web Conference 2024 (WWW-2024)*, pp. 1464–1473, 2024.

A VALID AND INVALID DEMONSTRATIONS

As shown in Table 2, the rows with a white background represent the Chain-of-Thought Reasoning methods, while the rows with a gray background represent the Program-of-Thought Computation methods. The demonstration problem is: *From the image, you can see that the shorter base of the trapezoid is 6ft, the longer base is 16ft, and the height is 12ft. Find the area of the trapezoid.* And the standard answer is: *132*.

Table 2: A pair of Valid and invalid demonstrations for different prompt methods in our experiment.

712	Prompting	Valid	Invalid		
713	CoT	In this image, since the trapezoid has a longer base (base1) with a	The image shows an equilateral triangle, which means all three sides		
714		length of 16 and a shorter base (base2) with a length of 6. The height,	are equal in length. Since one of the sides BC is given as 6 inches,		
715		which is the perpendicular distance between the two bases, is given as 12. To find the area of a trapezoid, you can use the formula: Area =	the other two sides (AB and AC) are also 6 inches each. To find the perimeter of the polygon we sum the lengths of all sides: Perimeter =		
716		1/2 * (base1 + base2) * height = 1/2 * (16 + 6) * 12=132. The area of	AB + BC + AC= $6+6+6=18$. The answer is 18.		
717		the trapezoid is 132. The answer is 132.			
718	LtM	Let's break down this problem: 1. What information can we obtain	A: Let's break down this problem:1.What information can we obtain		
719		from the graph? 2. How to use this information to calculate the area of this trapezoid? 1. The length of the longer base of this trapezoid	from the graph? 2. How to use this information to calculate the mea- sure of the perimeter of the regular polygon. 1. A regular triangle with		
720		is 16, the length of the shorter base is 6, and the height between the	one of the sides, AC, is given as 6 inches 2. Regular triangle, also		
721		longer and shorter base is 12. 2. According to the area formula of the transzoid the area of this transzoid is $0.5 * (16+6) * 12-132$. The	known as an equilateral triangle, which means all three sides are equal in length. Since one of the sides AC is given as 6 inches the other two		
723		answer is 132.	sides are also 6 inches each. We sum the lengths of all sides: Perimeter		
723			= AC + AB + BC = 6+6+6=18. The answer is 18.		
725	PS	Given: A trapezoidal with a longer base length of 16, a shorter base	Given: A regular triangle with one of the sides, AC, is given as 6		
726		length of 6, and a height of 12. Plan: We need to calculate the area of this teaconrid. Calculations the area formula of the teaconrid is 0.5	inches. Calculation: a regular triangle is also known as an equilateral		
727		* (longer base+shorter base) * height; the area of the trapezoid is 0.5	of the sides AC is given as 6 inches, the other two sides are also 6		
728		* (16+6) * 12=132. Answer: the area of the trapezoid is 132. The	inches each. We sum the lengths of all sides: Perimeter = $AC + AB$		
729		answer is 132.	+ BC=6+6+6=18. Answer: The perimeter of the triangle is 18. The answer is 18		
730	DD	According to the image the largely of the larger have of this tops	The impose shows an ampletant triangle, which means all three sides		
731	Kľ	zoid is 16, the length of the shorter base is 6, and the height be-	are equal in length. Since the first side AC is given as 6 inches, the		
732		tween the upper and shorter base is 12. According to the area formula	second and third sides (AB and BC) are also 6 inches each. To find the		
733		of the trapezoid, the area of the trapezoid is equal to 0.5 * (longer	perimeter of the polygon, we sum the lengths of all sides: Perimeter $AC + AB + BC = 6$ [the first side] if [the second side] if [the third		
734		(16[longer base]+6[shorter base]) * 12[height]=132. The answer is	side]=18. The answer is 18.		
735		132.			
736	PoT	base1 = 16	AC=6		
737		base2 = 6	AB=AC		
738		neight = 12 area = $0.5 * (base1+base2) * height$	BC=AC Perimeter=AC+AB+BC		
739		ans=area	ans=Perimeter		
740	PAL	# solution in Python:	# solution in Python:		
741		def solution():	def solution():		
742		Longer. base = 16 Shorter base = 6	Line_AC=6 Line_AB=Line_AC		
743		Height = 12	Line_BC=Line_AC		
744		Area = 0.5 * (Shorter_base+Longer_base) * Height	Perimeter=Line_AC+Line_AB+Line_BC		
745		return Area	return Perimeter		
746	SATLM	Shorter_base=6	Line_AC=6		
747		Height=12	Line_BC=Line_AC		
748		Area=Variable()	Perimeter=Variable()		
749		Area=0.5 * (Shorter_base+Longer_ base) * Height	Perimeter=Line_AC+Line_AB+Line_BC		
750		slove(result)	slove(result)		
751	MP	Mapping={Shorter_base:6,Longer_base:16,Height:12}	Mapping={Line_AC:6,Line_AB:6,Line_BC:6}		
752		# Algabraic answer	# Algabraic answer		
753		Answer=0.5 * (Shorter_base+Longer_ base) * Height	Perimeter=Line_AC+Line_AB+Line_BC		
754		# pymon code def solution(Shorter_base,Longer_ base,Height):	# python code def solution(Line_AC,Line_AB,Line_BC):		
755		return 0.5 * (Shorter_base+Longer_ base) * Height	Return Line_AC+Line_AB+Line_BC		

ANSWERING ACCURACY COMPARISON OF GPT-4V WITH ONE-SHOT В DEMONSTRATIONS

We evaluate the answering accuracy of GPT-4V through the prompting methods with one valid and one invalid demonstration (given in Table 2), respectively. As shown in Table 3, Compared to valid demonstration, prompting methods with invalid demonstration can sometimes improve the accuracy of GPT4-V answers, while at other times they can decrease them. For example, on the GeoQA+ dataset, the PS prompting method with a valid demonstration improved GPT-4V's average answer accuracy by 8.4% compared to an invalid demonstration, whereas it resulted in a 2.2% decrease on the GeoQA dataset. This indicates that the effectiveness of the demonstration is not a factor that affects the performance of the prompting method.

Table 3: Answering accuracy comparison of GPT-4V with one-shot valid and invalid demonstrations for solving geometric problems under different prompting settings on the five benchmark datasets.

770	Prompting	Setting	GEOS	Geometry3K	GeoQA	GeoQA+	PGPS9K
771	Average answering accuracy						
773	СоТ	Valid Reasoning	30.2	21.5	30.5	23.1	22.1
774		Invalid Reasoning	31.7	20.4	26.2	22.5	23.1
775	LtM	Valid Reasoning	28.6	21.2	30.5	20.1	20.5
776		Invalid Reasoning	30.7	23.5	26.5	25.0	20.5
777	PS	Valid Reasoning	37.8	28	32.3	30.5	23.3
778		Invalid Reasoning	36.2	30	34.5	22.1	25
779 780 781	RP	Valid Reasoning Invalid Reasoning	31.9 37.3	31.5 30.5	32.5 38.5	25.3 28.5	28.5 22.3
782	РоТ	Valid Computation	20.1	16.2	19.3	20.5	17.6
783		Invalid Computation	21.3	15.6	20.2	20.2	11.6
784	PAL	Valid Computation	20.8	16.4	20.3	19.3	11.9
785		Invalid Computation	19.1	17.5	21.5	18.8	11.1
786	SATLM	Valid Computation	36.2	22.4	28.7	23.1	22.5
787		Invalid Computation	35.3	21.9	30.9	23.3	22.7
788 789 790	MP	Valid Computation Invalid Computation	36.2 42.2	32.5 34.2	33.5 36.5	27.4 29.7	22.3 23.5
791		Self-co	onsistency	answering accı	ıracy		
792	СоТ	Valid Reasoning	38.3	36.7	45.3	48.3	26.7
793		Invalid Reasoning	37.3	36.5	46.2	47.4	27.1
794 795 796	LtM	Valid Reasoning Invalid Reasoning	41.7 42.3	40.9 40.1	43.3 44.7	34.6 39.7	33.3 31.6
797	PS	Valid Reasoning	43.7	31.4	37.8	35.8	25.9
798		Invalid Reasoning	44.5	31.9	35.3	35.3	29.1
799	RP	Valid Reasoning	45.3	26.7	40.3	40.7	33.3
800		Invalid Reasoning	45.5	28.1	41.4	40.4	34.7
801	РоТ	Valid Computation	23.9	19.3	26.2	25.7	20.4
802		Invalid Computation	23.6	18.7	26.9	24.1	19.8
803	PAL	Valid Computation	323.3	18.7	25.4	23.8	17.5
804		Invalid Computation	22.9	17.4	25.9	23.5	17.8
806	SATLM	Valid Computation	33.9	31.7	39.5	38.4	31.7
807		Invalid Computation	34.7	30.1	39.2	38.7	32.4
808	MP	Valid Computation	48.3	35.6	45.1	43.4	28.5
809		Invalid Computation	48.3	35.7	45.6	43.9	28.3

C ANSWERING ACCURACY COMPARISON OF GPT-4V WITH TWO-SHOT DEMONSTRATIONS

813 We evaluate the answering accuracy of GPT-4V through the prompting methods with two valid and 814 two invalid demonstrations, respectively. As shown in Table 4, Compared to valid demonstration, 815 prompting methods with invalid demonstration can sometimes improve the accuracy of GPT4-V 816 answers, while at other times they can decrease them. For example, on the Geometry3K dataset, the 817 CoT prompting method with a valid demonstration increased GPT-4V's average answer accuracy by 818 1.3% compared to an invalid demonstration, whereas it resulted in a 3.5% decrease on the PGPS9K 819 dataset. This indicates that the effectiveness of the demonstration is not a factor that affects the performance of the prompting method. 820

Table 4: Answering accuracy comparison of GPT-4V with two-shot valid and invalid demonstrations for solving geometric problems under different prompting settings on the five benchmark datasets.

824	Prompting	Setting	GEOS	Geometry3K	GeoQA	GeoQA+	PGPS9K	
825	Average answering accuracy							
827	СоТ	Valid Reasoning	33.5	22.5	32.5	31.3	21.5	
828		Invalid Reasoning	35.1	21.2	29.9	34.8	25.0	
829	LtM	Valid Reasoning	35.6	23.5	34.1	31.5	27.8	
830		Invalid Reasoning	36.8	30.7	36.5	31.7	25.9	
831	PS	Valid Reasoning	38.7	30.6	37.4	33.9	29.2	
832		Invalid Reasoning	36.3	27.8	36.9	37.5	30.6	
833 834 835	RP	Valid Reasoning Invalid Reasoning	40.7 41.8	31.5 30.1	41.5 44.9	40.3 38.7	30.5 27.8	
836	РоТ	Valid Computation	27.8	20.2	25.2	26.5	18.2	
837		Invalid Computation	29.1	18.4	24.5	20.2	21.4	
838	PAL	Valid Computation	25.2	18.9	22.7	23.5	21.9	
839		Invalid Computation	29.6	24.7	22.6	28.3	18.6	
840	SATLM	Valid Computation	31.5	27.0	32.7	31.9	24.5	
841		Invalid Computation	36.4	24.9	33.5	31.5	25.7	
842 843 844	MP	Valid Computation Invalid Computation	35.1 39.5	26.7 29.7	33.5 37.1	30.5 29.6	29.5 26.9	
845		Self-consistency answering accuracy						
846	СоТ	Valid Reasoning	49.8	30.2	45.2	44.5	38.2	
847		Invalid Reasoning	49.9	31.7	44.3	44. 4	39.9	
848	LtM	Valid Reasoning	50.3	35.9	45.8	46.8	38.7	
849		Invalid Reasoning	46.7	34.5	46.1	44.9	35.5	
851	PS	Valid Reasoning	50.7	38.2	47.4	46.2	38.4	
852		Invalid Reasoning	49.5	37.3	47.6	46.3	39.2	
853	RP	Valid Reasoning	53.1	38.9	49.5	47.7	39.4	
854		Invalid Reasoning	52.4	39.7	49.0	46.9	38.1	
855	РоТ	Valid Computation	32.5	30.0	30.8	31.1	27.1	
856		Invalid Computation	31.4	29.3	30.6	31.3	25.7	
857 858 850	PAL	Valid Computation Invalid Computation	31.8 30.6	29.2 28.8	30.6 31.9	32.4 32.7	29.9 29.3	
860	SATLM	Valid Computation	42.6	30.7	47.8	42.2	32.9	
861		Invalid Computation	41.7	31.2	48.6	41.5	33.7	
862	MP	Valid Computation	49.7	35.1	48.9	47.7	39.9	
863		Invalid Computation	49.5	34.5	48.6	48.1	40.4	

864 ANSWERING ACCURACY COMPARISON OF GPT-4V WITH THREE-SHOT D 865 DEMONSTRATIONS 866

867 We evaluate the answering accuracy of GPT-4V through the prompting methods with three valid and 868 three invalid demonstrations, respectively. As shown in Table 5, Compared to valid demonstration, 869 prompting methods with invalid demonstration can sometimes improve the accuracy of GPT4-V 870 answers, while at other times they can decrease them. For example, on the GeoQA+ dataset, the 871 SATLM prompting method with a valid demonstration increased GPT-4V's self-consistency answer 872 accuracy by 0.3% compared to an invalid demonstration, whereas it resulted in a 0.8% decrease on the GeoQA dataset. This indicates that the effectiveness of the demonstration is not a factor that 873 affects the performance of the prompting method. 874

875 Table 5: Answering accuracy comparison of GPT-4V with three-shot valid and invalid demonstra-876 tions for solving geometric problems under different prompting settings.

878	Prompting	Setting	GEOS	Geometry3K	GeoQA	GeoQA+	PGPS9K	
879	Average answering accuracy							
881	СоТ	Valid Reasoning	36.5	32.1	35.4	33.1	23.9	
882		Invalid Reasoning	36.2	27.5	41.1	29.2	23.5	
883	LtM	Valid Reasoning	38.4	28.5	34.5	29.5	24.5	
884		Invalid Reasoning	40.1	32.2	34.1	27.8	22.5	
885	PS	Valid Reasoning	39.2	31.7	33.8	32.5	27.5	
886		Invalid Reasoning	40.1	36.4	38.5	30.7	26.5	
887 888 889	RP	Valid Reasoning Invalid Reasoning	34.1 38.7	25.9 33.5	34.6 37.5	28.8 26.3	29.9 28.2	
890	РоТ	Valid Computation	28.7	22.1	26.8	26.5	20.3	
891		Invalid Computation	29.1	18.4	24.5	20.2	22.7	
892	PAL	Valid Computation	25.8	19.7	26.5	25.5	25.9	
893		Invalid Computation	29.2	25.5	25.5	27.5	20.5	
894	SATLM	Valid Computation	34.6	34.7	31.5	27.5	20.6	
895		Invalid Computation	36.8	27.5	34.2	29.1	27.5	
896 897 898	MP	Valid Computation Invalid Computation	32.9 34.6	34.6 30.5	35.7 39.1	28.4 31.4	24.5 24.1	
899		Self-consistency answering accuracy						
900	СоТ	Valid Reasoning	50.4	32.3	47.5	47.1	39.6	
901		Invalid Reasoning	51.1	32.4	47.1	46.7	39.7	
902	LtM	Valid Reasoning	51.4	35.1	48.9	47.1	40.1	
903		Invalid Reasoning	52.1	35.9	48.4	46.3	40.5	
904 905 906	PS	Valid Reasoning Invalid Reasoning	52.3 52.5	39.4 40.8	47.5 46.1	47.1 47.3	39.5 40.9	
907	RP	Valid Reasoning	55.7	40.9	50.1	49.3	49.1	
908		Invalid Reasoning	55.4	40.4	50.6	49.4	39.9	
909	РоТ	Valid Computation	35.7	30.4	36.8	34.3	29.5	
910		Invalid Computation	35.4	30.9	36.9	33.8	29.1	
911	PAL	Valid Computation	34.8	29.6	30.7	32.9	29.3	
912		Invalid Computation	34.9	30.1	30.1	32.7	30.2	
914	SATLM	Valid Computation	47.3	34.7	44.9	44.1	31.5	
915		Invalid Computation	47.2	35.9	45.7	43.8	31.9	
916	MP	Valid Computation	52.9	39.5	50.4	49.3	39.9	
917		Invalid Computation	53.6	39.1	50.9	49.8	38.7	





Figure 6: The distribution of the number of problems involving knowledge from different domains in two datasets, GeoQA and GeoQA+. DK_{-i} indicates that answering a geometry problem requires at least *i* domain knowledge.

F CASE ANALYSIS OF COMPUTATIONAL REQUIREMENTS

Table 6: Case analysis of computational requirements for different geometry problems.

Ш	Problem	Numerical values	Multi-dioit			
	text	image	answer	Tumerieur vurdes	intana uigit	
p_0	Find x.	x 47 57	0.429	$\frac{4}{7}, \frac{5}{7}, 0.429$	0	
p_1	Each pair of polygons is similar. Find the scale fac- tor from polygon ADCB to polygon PSRQ.		2.0	1.4, 0.7, 2.2, 3.2, 2	1	
p_2	Find the area of the shaded region. Assume that all polygons that appear to be regular are regular.	In	18.491	3, 18.491	2	
p_3	Find the measure of angle 1.	No. of the second se	112.0	34, 72, 112	3	
p_4	Find the area of the rhom- bus.		1200.0	20, 30, 1200	4	
p_5	Find the area of the circle.	222 n	38707.567	222, 38707.567	5	