
Assessing GPT4-V on Structured Reasoning Tasks

Mukul Singh¹ Jose Cambronero² Sumit Gulwani² Vu Le² Gust Verbruggen³

Abstract

Multi-modality promises to unlock further uses for large language models. Recently, the state-of-the-art language model GPT-4 was enhanced with vision capabilities. We carry out a prompting evaluation of GPT-4V and five other baselines on structured reasoning tasks, such as mathematical reasoning, visual data analysis, and code generation. We show that visual Chain-of-Thought, an extension of Chain-of-Thought to multi-modal LLMs, yields significant improvements over the vanilla model. We also present a categorized analysis of scenarios where these models perform well and where they struggle, highlighting challenges associated with coherent multimodal reasoning.

1. Introduction

Visual instruction tuning (Liu et al., 2023) allows language models to answer questions about images. By leveraging a dataset of captioned images, strong language models (like GPT-4) can be used to generate high-quality training data that spans conversation, detailed descriptions and reasoning about concepts found in these images. The visual instruction tuned model is then able to reason (to some extent) about the input image (Zhang et al., 2023c). OpenAI GPT-4V is the state-of-the-art language model GPT-4 (OpenAI, 2023) augmented with a vision component.

Example 1. The following image and part of the associated complex reasoning instruction are part of the LLaVA dataset (Liu et al., 2023). The image and **Q** are provided to the model, and it is trained to return **A**.



Q: What can be inferred about the relationship between these zebras?

A: The zebras are standing near each other and nuzzling their heads together behind the fence in their enclosure. This behavior suggests a bonding moment or a sign of affection between them.

In large language models, such as GPT-4, explicitly eliciting reasoning significantly improves performance. A popular method is Chain-of-Thought (CoT) where a literal instruction to “think step by step” causes the model to first generate

some thoughts, and then use these thoughts to generate a final answer (Wei et al., 2022). Similar to complex reasoning in LLaVA, multi-modal Chain-of-Thought (m-CoT) (Zhang et al., 2023c) leverages generating CoT data to fine-tune a multi-modal model to illicit reasoning about and answer scientific questions over graphics.

In this work, we explore the extent to which GPT-4V can *reason* in different domains, such as code and math, when part of the inputs are images or can be represented as images. More specifically, we evaluate GPT-4V in four domains.

1. Mathematical questions with visual context from the MathVista dataset (Lu et al., 2023).
2. Questions about data charts from the ChartQA dataset (Masry et al., 2022).
3. Abstraction and reasoning problems over objects in grids from the ARC dataset (Chollet, 2019; 2023).
4. Generating SQL from NL given a rendered table from the Spider dataset (Yu et al., 2018).

Each of these domains requires specific reasoning that might not be present in the training data.

Rather than fine-tuning a model, like m-CoT, we adapt the classical CoT approach to multi-modal language models, and show that it significantly outperforms the base (visual instruction tuned) model. Besides using the exact reasoning structure from m-CoT, we introduce an improved, three-step reasoning process to further improve performance, called *visual Chain-of-Thought* (v-CoT). In v-CoT, we instruct the model to (1) extract relevant information about the image, (2) use this relevant information to reason about the problem, and (3) state the final answer. Note that step (1) is not the same as a regular image captioning step, as it is also conditioned on the associated question along with the image.

Example 2. Figure 1 shows two examples where without reasoning, the model is not able find the correct answer.

Example 3. Figure 2 shows an example of reasoning with m-CoT and v-CoT. By first describing relevant information, the model is then able to reason about that information to obtain the right result.

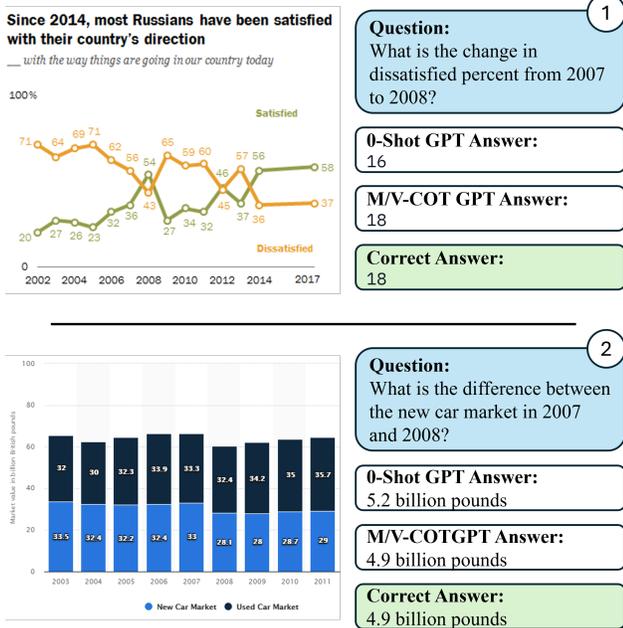


Figure 1. Without reasoning, the model does not find the correct answer. Both m-CoT and v-CoT elicit sufficient reasoning.

Our results show that for MathVista, ChartQA and Spider, GPT-4V outperforms baselines, including CoT and Program-of-Thought (PoT) prompting with GPT-4 over captions generated by InstructBlip (Dai et al., 2023). For ARC, we find that directly prompting over the vision component actually results in lower performance than captioning and CoT or PoT. We find that using v-CoT improves 1.5 – 9.3 percentage points over vanilla GPT-4V, and 1.1 – 4.7 percentage points compared with the m-CoT prompt.

In summary, we make the following contributions:

- A comparative evaluation of GPT-4V and (captioning + GPT-4) on structured reasoning tasks.
- Visual Chain-of-Thought (v-CoT) as an extension of CoT to multi-modal LLMs that first asks the model to extract relevant properties to reason over. v-CoT improves on vanilla GPT-4V by 1.5 – 9.3 percentage points.
- A thorough analysis of recurring patterns to understand GPT-4V’s performance.

2. Related Work

Prior work has explored combining multiple modalities using transformer-based models, including visual and text (Hu et al., 2023; Dai et al., 2023; Li et al., 2023; Lin et al.,

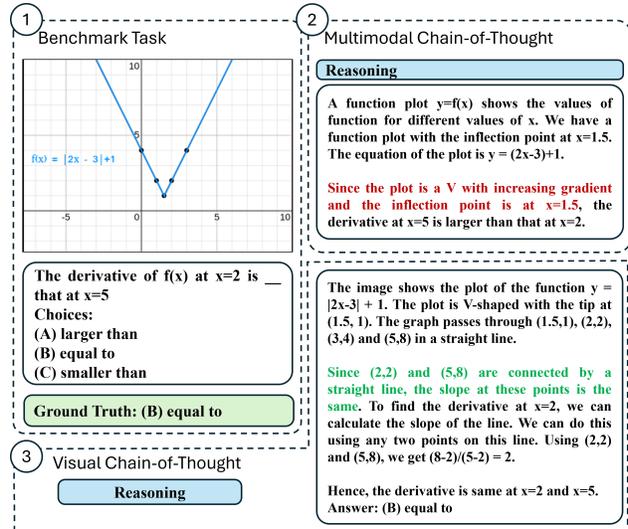


Figure 2. (1) An example task, (2) the output generated m-CoT, and (3) the output generated with v-CoT. Red and green text highlight incorrect and correct reasoning, respectively.

2023), speech and text (Zhang et al., 2023a), video and audio (Zhang et al., 2023b; Maaz et al., 2023), and the union of these modalities (Wu et al., 2023). In addition, there is a long history of models taking inputs in one modality (like image or text) and generating outputs in another (like text or images, respectively) (Hossain et al., 2019; Reed et al., 2016). For the image-text multimodal models, visual question-answering (VQA) (Antol et al., 2015; Agrawal et al., 2016) and conditioned captioning (Li et al., 2023; Dai et al., 2023) are the most popular tasks where the input is an image and optionally some text and the output is text conditioned on both the input image and text. In this work, we present an empirical evaluation of an existing state-of-the-art vision-text multimodal model (GPT-4V) on structured tasks like mathematical reasoning and code generation, and consider a setting similar to VQA where both images and text are used as input and the target output is text.

To evaluate these multi-modal models a large variety of benchmarks have been introduced, including VQA (Antol et al., 2015) and RefCOCO (Yu et al., 2016). In domains of structured reasoning, recent datasets include MathVista (Lu et al., 2023), which collects a large number of mathematical reasoning tasks, and ChartQA (Masry et al., 2022), which presents questions to be answered based on data analysis plots. Recently, MMMU (Yue et al., 2023) introduced college-level multi-modal tasks across a large range of domains including technical areas. Our evaluation also considers rendering existing structured inputs (tables) as images to perform text-to-SQL generation over them, using the popular Spider dataset (Yu et al., 2018).

Chain-of-Thought (Wei et al., 2022) reasoning traces were

used to fine-tune a multi-modal CoT model Zhang et al. (2023c). We take inspiration from m-CoT in our prompting experiments with GPT-4V. Leveraging a larger model (GPT-4), our visual guided prompting techniques are aimed at generalizing better to new modes of reasoning.

3. Evaluation Setup

First, we describe v-CoT and all baselines used in our evaluation. Second, we describe the benchmarks on which these approaches are evaluated.

3.1. Methods

We divide our baselines in three categories: multi-modal prompting strategies (m-CoT and v-CoT), directly using large instruction-tuned vision-text models, and captioning with instruction-tuned models plus reasoning over these captions with GPT-4.

3.1.1. M-COT AND V-COT

Figure 3 compares the corresponding prompts for v-CoT and m-CoT. The latter works as traditional CoT but over both text and image inputs. Specifically, the instruction to perform reasoning specifies that the image can be used to arrive at final answer in a step-by-step fashion.

The v-CoT prompt makes two small changes: (1) it instructs to describe relevant information and relevant image artifacts required to answer the question, and then (2) it asks to reason about this information to obtain the final answer. These artifacts act like predicates that can be used to reason over the image to solve the task, and can range from concrete values in the image (*the maximum value of the red line is 23*) to abstract concepts (*a right angled triangle*). Figure 2 shows an example where the model first identifies the points the line passes through even when this is not explicitly stated in the image.

3.1.2. INSTRUCTION TUNING

Sphinx (Lin et al., 2023) is a large multi-modal model (LMM) with multi-purpose visual instruction-following capabilities Sphinx has been trained on a variety of vision and language alignment tasks like object detection, visual question answering and region level captioning and achieves state-of-the-art performance on these.

Blip2 (Li et al., 2023) is a text and image pre-training technique that bootstraps vision-language pre-training from off-the-shelf frozen pre-trained image encoders and frozen large language models. The model combines existing vision and language models and pre-trains them on language and vision alignment tasks. Blip2 has been trained on various image-language tasks like conditional image captioning.

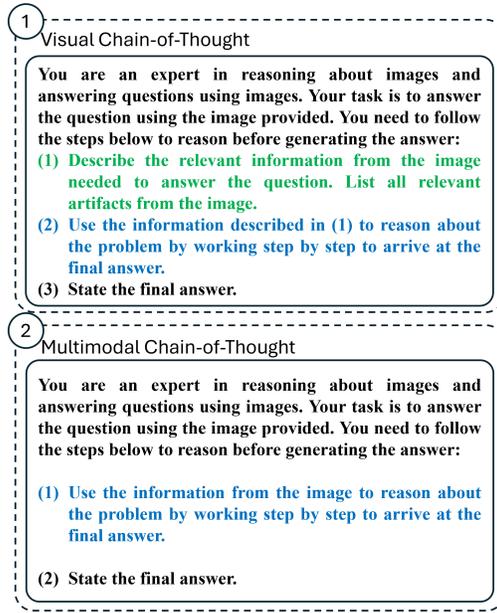


Figure 3. Prompt structure for v-CoT (1) and m-CoT (2) prompts. Blue highlights the shared m-CoT instruction and green highlights our extension.

InstructBlip (Dai et al., 2023) is an extension that applies instruction tuning to Blip-2 for language and vision tasks like visual question answering.

3.1.3. CAPTIONING + PROMPTING

InstructBlip + GPT4 with Chain-of-thought (CoT) (Wei et al., 2022) and **Program-of-thought (PoT)** (Chen et al., 2023) prompting of GPT-4. CoT and PoT have shown state-of-the-art performance in text and code generation tasks. To incorporate images, we first generate a caption for the image using InstructBlip, add the caption to the input prompt for GPT-4, and provide instructions for CoT or PoT.

3.2. Benchmarks

This section describes the benchmark tasks used in our evaluation. For all datasets, we perform exact match checks compared to the ground truth. We manually inspect failures to ensure they are not matching issues.

Table 1 summarizes our four benchmarks, and Figure 4 presents an example task from benchmark. To mitigate computational costs, we sample 20% tasks uniformly at random from each dataset.

3.2.1. MATHEMATICAL REASONING

We evaluate mathematical reasoning given a visual context with the recently introduced MathVista dataset (Lu et al., 2023). These problems require reasoning over images which

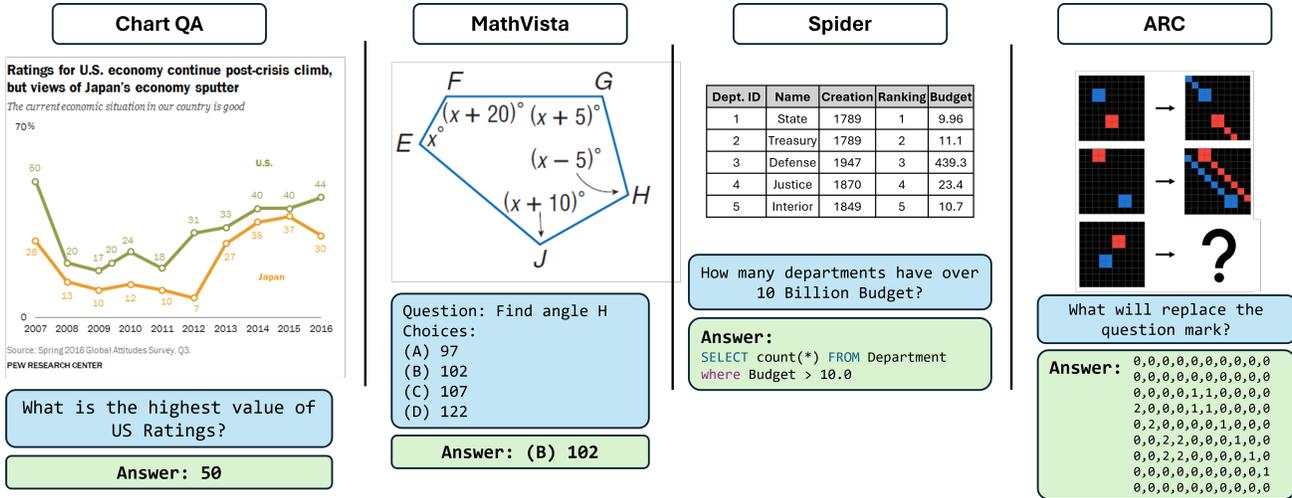


Figure 4. Sample tasks from each benchmark dataset. We show the image and the associated text prompt for the dataset along with the correct answer for the task. For Chart QA and MathVista the answer is a choice or numeric value; For Spider the answer is the correct SQL query; For ARC the output is the correct pixel matrix (0 → black; 1 → red; 2 → blue).

Dataset	Tasks	Image	Task
ChartQA (test)	302	Chart	VQA
MathVista (testmini)	200	Diagram	VQA
ARC (test)	80	Blocks	Pattern
Spider (test)	206	Table	Code generation

Table 1. Tasks from our evaluation benchmarks, we present the total number of tasks, image description, and task description. We sample 20% from each dataset to mitigate computational costs.

may contain diagrams, tables or other images. Each task consists of an image and question pair. Questions in MathVista can be multiple-choice, free-form with single value as result, or free-form with a list as a result. We sampled tasks from the testmini split of MathVista, which resulted in multiple choice and single-value free-form questions¹.

3.2.2. VISUAL DATA ANALYSIS

We evaluate the ability to answer questions based on charted data—an important skill for visual data analysis—using the ChartQA dataset (Masry et al., 2022). Each task consists of a plot as a rendered image and a question about this plot. These questions are short and objective, with 2-3 word or numeric value answers. An example question is “How many crimes were committed in 2020?” The question is passed in the prompt, preceded by the following instruction:

Answer the question using the image. Only give the exact answer in 2-3 words or as a numeric value.

¹testmini has only 2 free-form list questions

3.2.3. VISUAL ABSTRACTION AND EXTRAPOLATION

We evaluate the ability of multi-modal models to solve program-synthesis-like tasks over grids that require abstraction and extrapolation by using the ARC dataset (Chollet, 2023). Each ARC task consists of a set of examples, where an example consists of an input and output grid, and a new input—the model should predict the associated output grid. Grids are represented as comma-separated-value (CSV) table of integers. We use the following instruction

Generate the transformed representation that will replace the question mark by looking at the example figure transformation. Generate as a <gridsize> Grid where 0 denotes black, 1 denotes red and 2 denotes blue.

where we replace <gridsize> with the corresponding dimensions specified in the individual benchmark task.

3.2.4. CODE GENERATION

We evaluate NL-to-SQL generation in a multi-modal setting using Spider dataset (Yu et al., 2018). Each task input consists of a relational database and natural language question pair, and the task output is the SQL code needed to answer that question. To turn Spider tasks into multi-modal tasks, we consider only tasks over a single table and we render the first 50 rows into an image using matplotlib.pyplot.table (Hunter, 2007). We use the following instruction:

Generate SQL query for the given user question. The relevant table is shown in the image and the metadata is included.

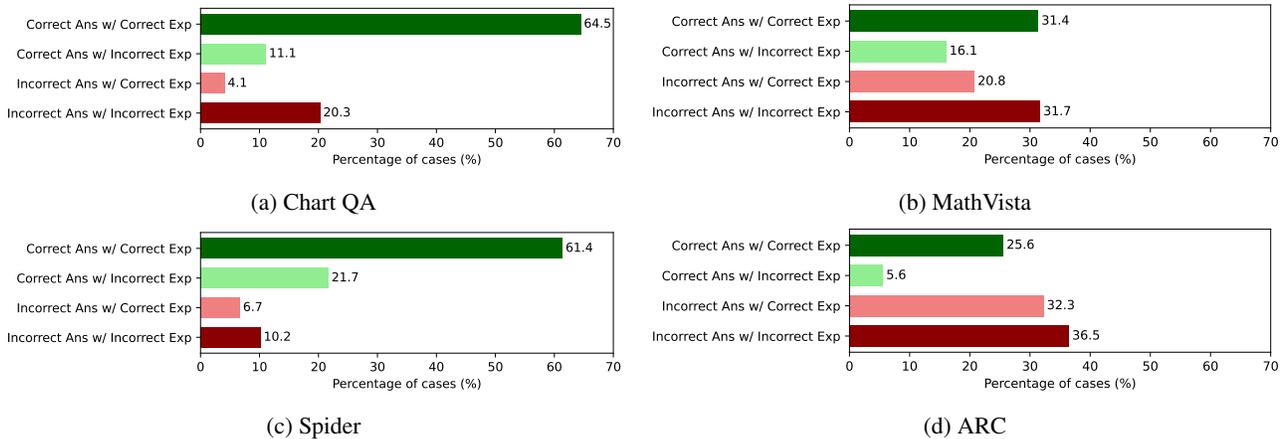


Figure 5. Manual analysis of GPT-4V + VCoT on our sampled tasks. We manually annotate the reasoning and the final answer separately for all benchmark dataset and present the analysis. For Chart QA and Spider, we find that in 72 – 85 % cases the model generates the right explanation and is able to generate the correct answer from these. For MathVista and ARC that require more complex reasoning, we find that the model struggles to generate the right reasoning with only 51 – 58 % reasoning being correct.

Approach	MathVista	ChartQA	ARC	Spider
Sphinx	33.2	63.4	27.2	-
InstructBlip	24.7	54.5	23.4	-
Blip2	23.1	31.4	22.6	-
InstructBlip + CoT GPT-4	30.4	55.3	47.3	83.5
InstructBlip + PoT GPT-4	30.8	56.4	51.2	82.7
GPT-4	-	-	-	81.8
GPT-4V	47.5	75.6	31.2	84.3
GPT-4V + v-CoT	49.1	79.2	40.5	86.2

Table 2. Results over 20% of tasks sampled from each benchmark uniformly at random. InstructBlip + GPT-4 indicates captioning by InstructBlip and reasoning over this by text GPT-4. CoT represents Chain of Thought; PoT represents Program of Thought. GPT-4 does not use a vision component and is only evaluated on problems that do not depend on an image.

Metadata: <Metadata>
 Question: <Question>

where we replace <MetaData> with the table schema (column names and types) and <Question> with the user question. We use the schema encoding from FormT5 (Singh et al., 2023). Note that we only evaluate GPT-4-based models on this problem, as other baseline methods have not been trained to generate code.

4. Results

Table 2 presents a summary of our results across benchmarks and baselines. We find that for both MathVista and ChartQA, GPT-4V outperforms other approaches. Interestingly, for the ARC dataset, we observe that using the multi-modal version of GPT-4 actually produces worse results, compared to generating textual descriptions and then applying GPT-4

plus CoT or PoT (the latter of which improves results most). Overall, the ARC dataset represents a challenging task. On the Spider dataset, we find that including image improves performance by 0.8 percentage points, while adding v-CoT improves by 1.9 percentage points.

We manually annotated GPT-4V results for whether the reasoning provided was correct or not and whether the answer was correct or not. Figure 5 presents our findings. Across all datasets, we observe that there are tasks that GPT-4V correctly answers, despite having incorrect reasoning. Similarly, we find that despite producing correct reasoning, the model can still produce an incorrect answer. This issue is particularly frequent in both the MathVista and ARC tasks.

4.1. Discussion of Patterns

We discuss recurring patterns we observed in our manual analysis of results.

4.1.1. MATHEMATICAL REASONING

We identify the following recurring patterns in the MathVista failures.

- Incorrect computation over derived values. Figure 7b shows an example where the model identifies correct values and constraints for a trigonometric problem, but fails to compute the final answer.
- Inferring numerical labels. When the model has to infer numerical labels, like aggregates in bar charts, the model often picks wrong values and fails to answer, despite having correct reasoning. This is shown in Figure 7c.

Assessing GPT4-V on Structured Reasoning Tasks

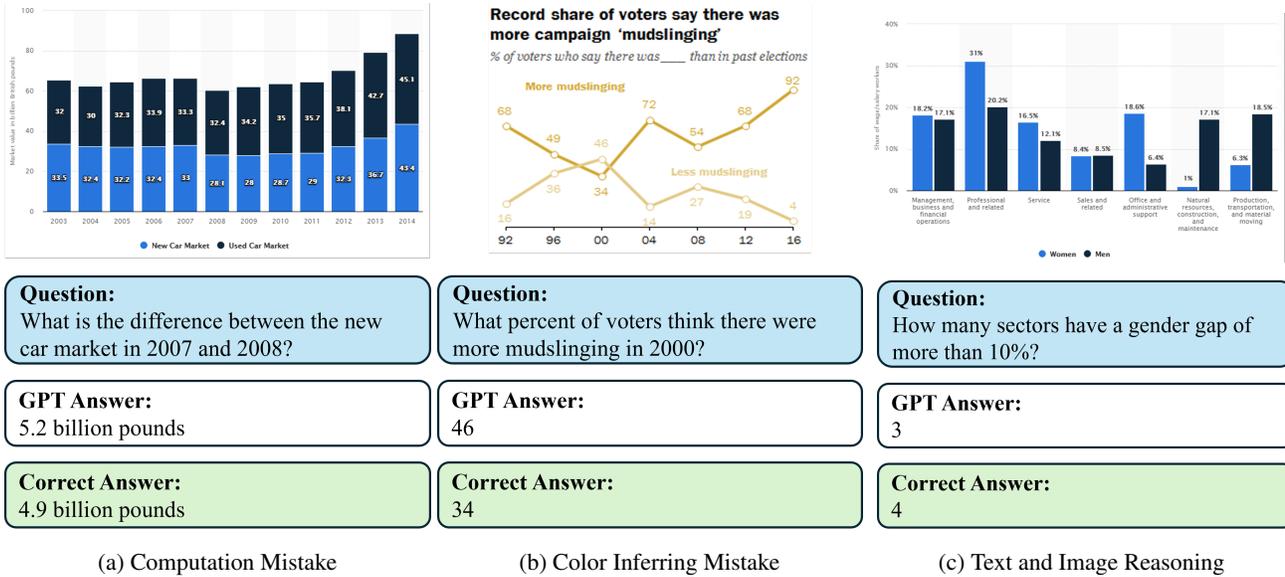


Figure 6. Observed patterns in GPT-4V behavior for ChartQA tasks. (a) shows an example where the model generates the correct reasoning, i.e. to subtract new car market share of 2008 from 2007, but makes a computation error in the final answer. (b) shows an example where the model generates the correct reasoning but applies it to the light yellow line instead of the dark yellow. (c) shows an example where the model needs to reason on both text and images together. V-CoT is able to fix this example.

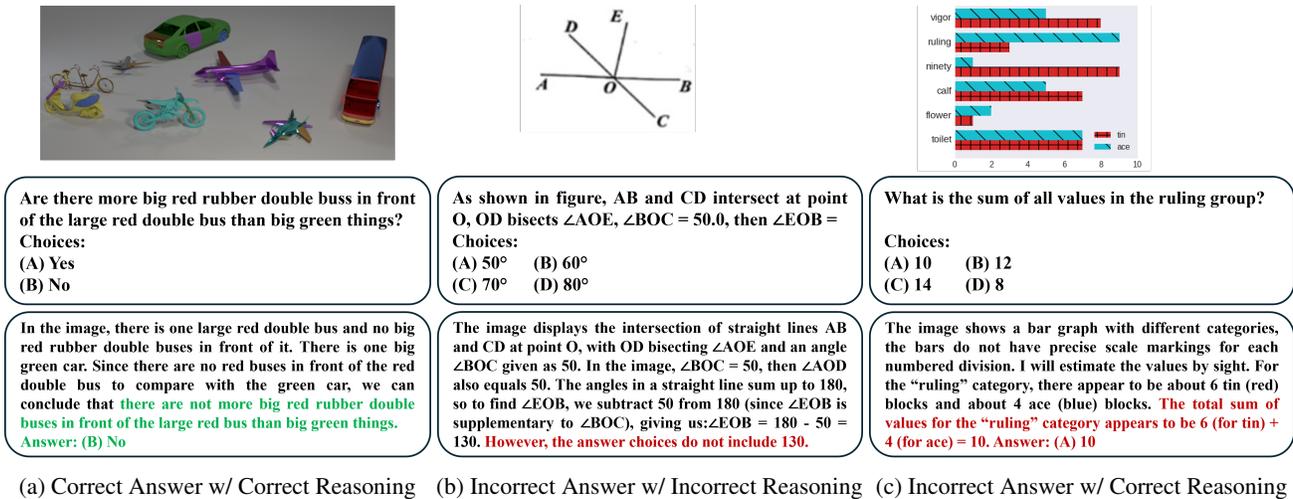


Figure 7. Observed patterns in GPT-4V + V-CoT responses to MathVista tasks. (a) shows an example where both the reasoning and answer are correct; (b) shows an example where both the reasoning and answer are incorrect as the model misses to subtract angle DOE; (c) shows an example where the answer is incorrect but the reasoning is correct, as the model correctly figures out the correct bar but cannot generate the correct value from scale.

4.1.2. VISUAL DATA ANALYSIS

We find that in ChartQA, multiple failures can be attributed to a few emergent behaviors, which are highlighted below. These patterns are also illustrated in Figure 6.

- Incorrect computation over chart-derived values. For example, take the question “What is the growth from 2010 to 2011?” over the chart in Figure 6a. Despite

identifying each of these values correctly, it fails to compute the growth amount by subtracting the value of 2007 from that of 2008.

- Linking colors/visual elements. For example, when asked a question about a plot with two lines of varying shade of yellow, the model flips their labeling and answers with the wrong value (see Figure 6b).

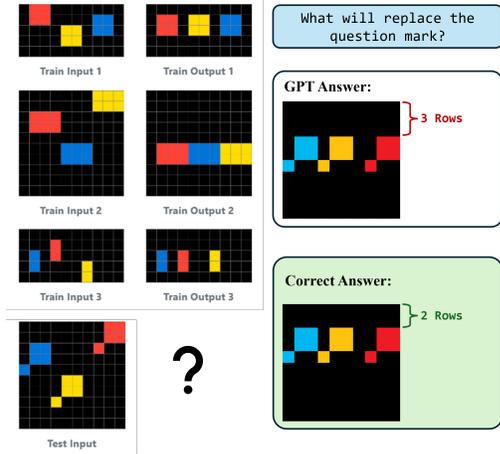


Figure 8. In ARC, GPT-4V can still struggle in identifying/placing blocks in appropriate locations (like offsets). The figure shows that GPT-4V correctly identifies that it has to align the blocks to the level of the blue block, but it is offset by one block.

- Computations that require extracting values from both image and the query. For example, extracting a rate from the text query and using this to identify differences in a plot (see Figure 6c).

4.1.3. ABSTRACTION AND EXTRAPOLATION

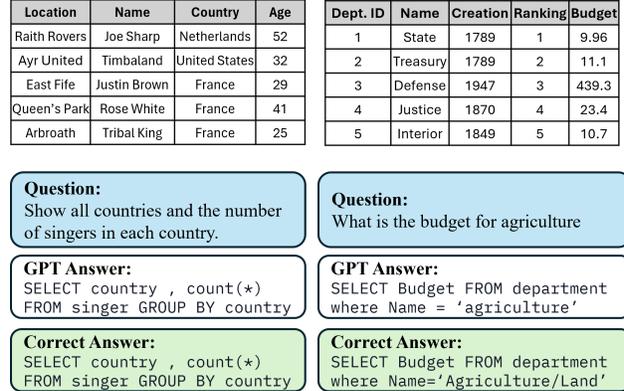
In the ARC dataset, all methods struggle to generate the correct output since it requires complex reasoning over visual patterns. We find the following to be a common pattern.

- Incorrect grid perception. Despite using v-CoT, we find cases where GPT-4V does not correctly identify block positions. For example, v-CoT may place blocks at the wrong vertical offset (see Figure 8).

4.1.4. CODE GENERATION

On SQL generation for Spider tasks, we make the following observations:

- Structural understanding. The rendered table provides improvements over the text representation of the same table for tasks that require reasoning over the structure of a table (pivot, group, transposing). An example is shown in Figure 9a.
- Token efficiency of images. Considering tables as images consumed fewer tokens than their corresponding text description produced by InstructBlip.
- Reasoning over cell values. Tasks that require reasoning over cell values, such as filtering or splitting, benefit from a textual representation of tables or from the image paired with v-CoT (see Figure 9b).



(a) Success Case

(b) Failure Case

Figure 9. We show success and failure cases of adding tables as images. (a) Adding the table image can help for queries that requiring overall table structure (like grouping) or extracting particular values (like filtering) from table rendering. (b) For queries that need constant values, GPT-4V struggles to get it from the image.

5. Conclusion

We present the first evaluation of the state-of-the-art multi-modal LLM GPT-4V on structured reasoning tasks across various domains. We show that our visual Chain-of-Thought (v-CoT) prompt improves performance by first instructing the model to analyse the image, conditioned on the task at hand, and then use this analysis to reason towards a final result. Our experiments compare GPT-4V (with and without v-CoT) to multiple existing multi-modal models. We find that v-CoT outperforms other baselines in three datasets, but that it still struggles with the ARC dataset, which requires abstraction and extrapolation.

6. Limitations

First, we consider a subset of structured reasoning tasks, performance on other tasks may vary. Similarly, we sample 20% of tasks per dataset to mitigate computational costs—while this sampling was done uniformly at random, it is possible that performance would change for other tasks. While manually annotating results, we also identified that there are tasks in the various datasets (particularly ChartQA and MathVista) that are ambiguous or may have incorrect ground truth labels. We take all labeling as given. All our tasks are framed in English (images with English text, tasks with English text). Evaluating GPT-4V and other models on multi-modal, multilingual tasks remains future work.

References

Agrawal, A., Lu, J., Antol, S., Mitchell, M., Zitnick, C. L., Batra, D., and Parikh, D. Vqa: Visual question answering,

- 2016.
- Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C. L., and Parikh, D. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pp. 2425–2433, 2015.
- Chen, W., Ma, X., Wang, X., and Cohen, W. W. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=YfZ4ZPt8zd>.
- Chollet, F. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- Chollet, F. The abstraction and reasoning corpus (arc), 2023. URL <https://github.com/fchollet/ARC>. last accessed: 15-11-2023.
- Dai, W., Li, J., Li, D., Tiong, A. M. H., Zhao, J., Wang, W., Li, B., Fung, P., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- Hossain, M. Z., Sohel, F., Shiratuddin, M. F., and Laga, H. A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):1–36, 2019.
- Hu, W., Xu, Y., Li, Y., Li, W., Chen, Z., and Tu, Z. Bliva: A simple multimodal llm for better handling of text-rich visual questions. *arXiv preprint arXiv:2308.09936*, 2023.
- Hunter, J. D. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007. doi: 10.1109/MCSE.2007.55.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.
- Lin, Z., Liu, C., Zhang, R., Gao, P., Qiu, L., Xiao, H., Qiu, H., Lin, C., Shao, W., Chen, K., et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv preprint arXiv:2311.07575*, 2023.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *NeurIPS*, 2023.
- Lu, P., Bansal, H., Xia, T., Liu, J., Li, C., Hajjishirzi, H., Cheng, H., Chang, K.-W., Galley, M., and Gao, J. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models, 2023.
- Maaz, M., Rasheed, H., Khan, S., and Khan, F. S. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- Masry, A., Do, X. L., Tan, J. Q., Joty, S., and Hoque, E. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.177. URL <https://aclanthology.org/2022.findings-acl.177>.
- OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023. URL <https://api.semanticscholar.org/CorpusID:257532815>.
- Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., and Lee, H. Generative adversarial text to image synthesis. In *International conference on machine learning*, pp. 1060–1069. PMLR, 2016.
- Singh, M., Cambronero, J., Gulwani, S., Le, V., Negreanu, C., Nouri, E., Raza, M., and Verbruggen, G. Format5: Abstention and examples for conditional table formatting with natural language, 2023.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., brian ichter, Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain of thought prompting elicits reasoning in large language models. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_VjQlMeSB_J.
- Wu, S., Fei, H., Qu, L., Ji, W., and Chua, T.-S. Nextgpt: Any-to-any multimodal llm. *arXiv preprint arXiv:2309.05519*, 2023.
- Yu, L., Poirson, P., Yang, S., Berg, A. C., and Berg, T. L. Modeling context in referring expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pp. 69–85. Springer, 2016.
- Yu, T., Zhang, R., Yang, K., Yasunaga, M., Wang, D., Li, Z., Ma, J., Li, I., Yao, Q., Roman, S., Zhang, Z., and Radev, D. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task. In Riloff, E., Chiang, D., Hockenmaier, J.,

and Tsujii, J. (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3911–3921, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1425. URL <https://aclanthology.org/D18-1425>.

Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., Wei, C., Yu, B., Yuan, R., Sun, R., Yin, M., Zheng, B., Yang, Z., Liu, Y., Huang, W., Sun, H., Su, Y., and Chen, W. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.

Zhang, D., Li, S., Zhang, X., Zhan, J., Wang, P., Zhou, Y., and Qiu, X. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. *arXiv preprint arXiv:2305.11000*, 2023a.

Zhang, H., Li, X., and Bing, L. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023b.

Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G., and Smola, A. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023c.