

MM-PhyQA: Multimodal Physics Question-Answering With Multi-Image CoT Prompting

Avinash Anand¹[0009-0003-2479-034], Janak Kapuriya¹[0009-0007-9562-7672], Apoorv Singh¹[0009-0006-2589-8293], Jay Saraf¹[0009-0003-3177-149X], Naman Lal¹[0009-0008-2914-5509], Astha Verma¹[0000-0003-3615-5373], Rushali Gupta¹[0009-0006-1399-8262], and Rajiv Shah¹[0000-0003-1028-9373]

Indraprastha Institute of Information Technology, Delhi, India

{avinasha, kapuriya22032, apoorv17027, jay20438, asthav, rajivrtn}@iiitd.ac.in

Abstract. While Large Language Models (LLMs) can achieve human-level performance in various tasks, they continue to face challenges when it comes to effectively tackling multi-step physics reasoning tasks. To identify the shortcomings of existing models and facilitate further research in this area, we curated a novel dataset, **MM-PhyQA**, which comprises well-constructed, high school-level multimodal physics problems. By evaluating the performance of contemporary LLMs that are publicly available, both with and without the incorporation of multimodal elements in these problems, we aim to shed light on their capabilities. For generating answers for questions consisting of multimodal input (in this case, images and text) we employed Zero-shot prediction using GPT-4 and utilized LLaVA (LLaVA and LLaVA-1.5), the latter of which were fine-tuned on our dataset. For evaluating the performance of LLMs consisting solely of textual input, we tested the performance of the base and fine-tuned versions of the Mistral-7B and LLaMA2-7b models. We also showcased the performance of the novel **Multi-Image Chain-of-Thought (MI-CoT)** Prompting technique, which when used to train **LLaVA-1.5 13b** yielded the best results when tested on our dataset, with superior scores in most metrics and the highest accuracy of 71.65% on the test set.

Keywords: Large Language Models · Large Multimodal Models · Prompt Engineering · Chain-of-Thought

1 Introduction

Recent advances in Large Multimodal Models (LMMs) show impressive capabilities in handling multiple modalities, excelling in tasks like zero-shot generalization, visual reasoning, and instruction-following. Models like LLaMA-2 [1] and Mistral-7b [2] have displayed decent performance on famous textual mainstream question-answering benchmarks. SciPhyRAG [3] used retrieval augmentation to solve physics questions. However, the challenge of effectively handling queries combining textual and visual components persists, especially in subjects like Math and Physics, a problem that is exemplified by state-of-the-art models like GPT-4 [4] being proprietary. Fine-tuning general-purpose LLMs to perform well at a singular task has been effective in a variety

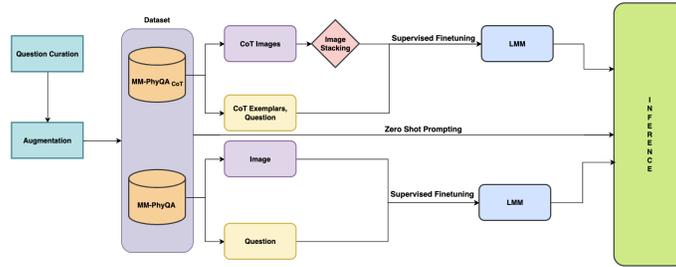


Fig. 1: Schematic Pipeline of Multimodal Question Answering

of complex scenarios [5,6]. Hence, developing open-source domain-specific chatbots with multimodal capabilities is promising. These chatbots can empower students with interactive question sessions, providing instant clarifications and guidance, and revolutionizing exam preparation.

To evaluate the capabilities of Large Multimodal Models (LMMs) for question-answering we have created a novel multimodal multiple-choice high school physics question-answering dataset. Physics questions require a good understanding of the underlying concepts and construction of steps with reasoning to reach the correct solution, hence not solvable by simply memorizing certain facts. High-school physics numerical questions are often accompanied by diagrams, which adds additional complexity that models should be able to interpret and understand for effective problem-solving, therefore acting as a valuable benchmark for evaluating the performance of LMMs. Given the dearth of multimodal physics datasets containing complex, high-quality questions, our dataset facilitates the study performance of LMMs and LLMs in a challenging setting.

Introduction of techniques like Chain-of-Thought (CoT) Prompting [7] has further enhanced the performance of LLMs, and subsequent experiments using the technique in a multimodal context [8,9] have been fruitful. CoT-Prompting involves providing the necessary prompts to a model to steer it toward the correct solution. It is analogous to how humans go about solving a problem, wherein we try to think of the intermediate steps that build logically toward the final answer. However, the prospect of incorporating images and figures with the prompt exemplars is yet to be explored by contemporary literature.

In this paper, we do a quantitative analysis regarding the effect of utilizing a modality other than text and the difference in the performance of LLMs and LMMs between using them out of the box (Zero Shot Prompting) and fine-tuning them for a specific purpose. We also examine the effects of using Chain-of-Thought Prompting in a multimodal setting, for which we came up with a novel method to incorporate multiple images during the CoT prompting process.

Hence, the contributions of this paper are threefold. Firstly we introduce a novel multimodal dataset, MM-PhyQA, containing challenging physics questions. We also generate its CoT-Prompting variant, providing exemplar questions during the training process. Secondly, we analyze the effects of using an additional modality other than

text, the effects of utilizing techniques like CoT Prompting, and the performance gain witnessed by fine-tuning LLMs and LMMs for a specific purpose, particularly for a task like answering physics questions. Finally, we introduce an approach, **Multi-Image Chain-Of-Thought (MI-CoT)** for employing multiple images during CoT-Prompting that is novel, to the best of our knowledge.

2 Related Works

2.1 Available Datasets

Numerous educational datasets are available for math and science. GSM8k [10] offers 8500 grade school math problems, while JEEBench [11] provides 450 questions from JEE advanced exams. SciQ [12] contains 13,697 science questions, and SciBench [13] offers college-level scientific problems. MMLU [14] is a multitask test dataset with 15908 samples, and C-Eval [15] includes multiple-choice questions in Chinese across 52 disciplines.

In the realm of multimodal datasets, GeoQA [16] offers middle school geometric questions with images and text, while TQA [17] provides middle school science questions in a similar format. ChartQA [18] is a chart-based reasoning dataset, and MMQA [19] consists of questions with images, text, and tables. ScienceQA [8] is a diverse multimodal dataset with 21208 science questions spanning various topics but lacks challenging high school-level questions.

2.2 Large Multimodal Models and Chain-of-Thought

Large language models' extension into multi-modal versions has led to significant attention and successful applications. GPT4-V [20] and PaLM-E [21] are state-of-the-art multimodal models, with PaLM-E directly incorporating visual features for enhanced performance. LLaVA [22,23] is recognized for its versatility in handling various multimodal tasks, utilizing a CLIP [24] encoder with Vicuna for vision-language understanding. Shikra [25] excels in Visual Question Answering (VQA) and image-captioning tasks, particularly in multimodal conversation scenarios. Kosmos-2 [26] demonstrates strong performance across diverse multimodal tasks, including grounding, referring, learning within context, and generation.

The Chain-of-Thought paradigm has transformed how large language models process reasoning, significantly improving NLP tasks. It has evolved from vanilla CoT to more complex structures like Tree-of-Thoughts [27] and Graph-of-Thoughts [28]. Despite these advancements, the shift towards multimodal reasoning led by multimodal CoT [9], has limitations due to reliance on multiple question-answer chains from a single image during training. To overcome this, we propose the Multi-Image Chain-Of-Thought (MI-CoT) technique, ensuring each question-answer pair used in training is associated with a distinct image, enhancing diversity and robustness.

3 Novel Dataset

There is a lack of multimodal datasets that comprise physics questions and are catered to high school students. While there are a few datasets available that consist of questions

Question: A ball of mass (m) 0.5 kg is attached to the end of a string having length (L) 0.5m. The ball is rotated on a horizontal circular path about vertical axis. The maximum tension that string can bear is 324 N. The maximum possible value of angular velocity of ball (in rad/s) is:

Options:

- a) 9
- b) 18
- c) 27
- d) 36

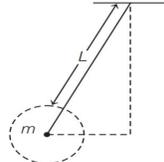


Figure 1: Image

Answer: (d)

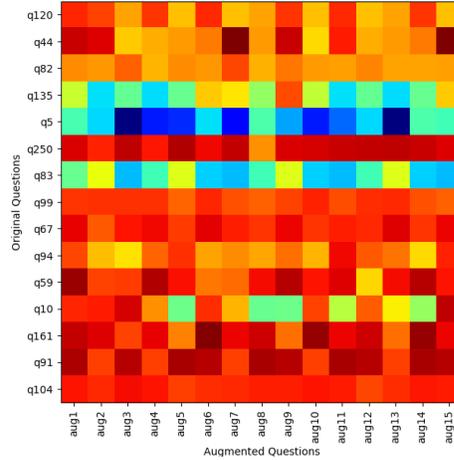
Explanation: From the figure, $T \sin \theta = mL \sin \theta \omega^2$.

$$324 = 0.5 \times 0.5 \times \omega^2$$

$$\omega^2 = \frac{324}{0.5 \times 0.5}$$

$$\omega = \sqrt{\frac{324}{0.5 \times 0.5}} = \frac{18}{0.5} = 36 \text{ rad/s}$$

(a) Sample question of MMPhy-QA dataset



(b) Heatmap of text similarity between 15 randomly sampled original and augmented questions

Fig. 2: MMPhy-QA Dataset questions

at a high school level, the quality of the questions does not belong to the highest of standards. We curated a novel MM-PhyQA Dataset from publicly available resources. The resources are geared toward individuals who prepare for competitive exams throughout India, ensuring a higher difficulty level than that of an average high school physics question.

3.1 Original Dataset Creation

Around 300 questions were manually created. As shown in Figure 2a, each question consists of a question, four options, the correct answer to the question, and an explanation that shows the reasoning by giving steps to approach the correct answer to select the correct answer.

3.2 Data Augmentation Procedure

For augmenting the data ChatGPT [29] was given a prompt to create other variations of the text while ensuring that the meaning remained the same, bringing the total count of

the questions in the dataset to 4500. Figure 2b shows the heatmap of the cosine similarity scores of the augmented questions w.r.t the original one for some of the questions. The questions were altered in two ways:

- **Numerical Value Variation:** During augmentation, numerical values in the original questions are adjusted to diversify the solutions, ensuring the model’s impartiality. Python functions were developed for each question to get the correct answers after changing the values.
- **Structural Variation:** To avoid pattern memorization, the questions’ structure was intentionally altered by rephrasing with ChatGPT and sometimes manual adjustments. Options were kept the same but randomly rearranged.

Initially, attempts to rephrase the entire query sometimes failed to properly shuffle the questions. Manual adjustments were made to correct these errors. While including the entire query didn’t consistently result in a rephrased version, prompting ChatGPT to generate separate variations for the question and explanation improved results. However, some questions still required manual rephrasing, involving adjustments to the question, explanation, options, and correct answer.

3.3 Chain of Thought Variant

To facilitate the model to generate better reasoning, two questions were added corresponding to each question. These questions were based on the same topic and care was taken that similar concepts were utilized as seen in Figure 2a. All three questions consist of figures.

Table 1: Topics and subtopics in the MM-PhyQA dataset

| Topic | Subtopics |
|---|---|
| Kinematics | Velocity-Time, Acceleration, Rotational Motion, Gravitation, Motion in a Straight Line, Motion in a Plane, Periodic Motion, Wave Motion. |
| Mechanics | Law of Motion, Work, Power, Force, Law of Motion |
| Electrostatics and Current Electricity | Current, Voltage, Resistance, Electric Field, Ohm’s Law, Kirchoff’s Laws, and Their Applications, Series and Parallel Combinations of Resistors |
| Thermodynamics | Laws of Thermodynamics, Thermal Equilibrium, Heat Transfer, Temperature, Reversible and Irreversible Processes, Kinetic Theory of Gases. |
| Optics | Reflection, Mirrors, Lenses, Wave Optics, Magnification. |
| Magnetism | Magnetic Field, Hysteresis, Permeability, Electromagnets. |
| Electronic Devices | Semiconductors, Logic Gates, Diode. |
| Atoms | Nuclei, Isotopes. |

3.4 MM-PhyQA Dataset Topics

The dataset consists of topics that are present in high school physics curricula throughout India. The topics and the corresponding subtopics are listed in Table 1.

| question | label | image |
|---|---|--|
| <p>Q30_1- The velocity - time graph of a particle moving in a straight line is shown in figure. The mass of the particle is 2kg. Work done by all the forces acting on the particle in time interval between $t = 0$ to $t = 10$ s is. , select the correct option from - a)300 J ; b)-300 J ; c) 400 J ; d)-400 J and give me the reason behind the selected option?</p> <p>the correct option is a and the reason is From work-energy theorem, $W = \Delta KE = K_f - K_i = \frac{1}{2} m (v_f^2 - v_i^2) = \frac{1}{2} \times 2 [(20)^2 - (10)^2] = 300 \text{ J}$</p> | <p>b- Initial Velocity of particle, $v_i = 20 \text{ ms}^{-1}$ Final velocity of the particle, $v_f = 0$ According to work-energy theorem, $W_{(net)} = \Delta KE = K_f - K_i$</p> | <p>q30_1.png, q30_2.png, q30.jpg</p> |
| <p>Q30_2- The v-t graph of a particle moving along the x-axis is shown in the figure. The mass of the particle is 4 kg. The work done by all the forces acting on the particle between $t = 3$ s to $t = 6$ s is. , select the correct option from - a)12 J ; b) 24 J ; c) 8 J ; d) 32 J and give me the reason behind the selected option?</p> <p>the correct option is b and the reason is By the equation of a line $y = mx + c$, we get, $v = -2t + 8$ At $t = 3$ s, $v = 2 \text{ m/s}$, $KE_i = 8 \text{ J}$ At $t = 6$ s, $v = -4 \text{ m/s}$, $KE_f = 32 \text{ J}$ $W = KE_f - KE_i = 24 \text{ J}$</p> | <p>$m(v_f^2 - v_i^2) = \frac{1}{2} \times 4 [(0)^2 - 20^2]$</p> | |
| <p>Q3- Velocity-time graph of a particle of mass 2 kg moving in a straight line is as shown in figure. Work done by all forces on the particle is., select the correct option from - a)400 J ; b)-400 J ; c) -200 J ; d)200 J</p> | | |

Fig. 3: Multi-Image Chain of thought (MI-CoT) Prompted text provided as input to LMMs during training. The main question to be answered is preceded by two exemplars, with the three questions separated by a delimiter. The image is a sequence of three comma-separated file names and the label is the ground truth

4 Methodology

Figure 1 shows the pipeline that was utilized for data processing, input processing, and output generation. Each element in the dataset consists of the question ID, the question, the label consisting of the corresponding answer and the reasoning, and the image filename. A function was used to convert each element to a prompt which can be fed to the model for generating the answer. For the Chain of Thought variant of the dataset, the structure was modified. As shown in Figure 3, the question was preceded by two similar questions with their correct answers and reasoning. All the three questions were separated by a delimiter consisting of hyphens. The filenames of the three images were stored in a comma-separated fashion.

4.1 Multi-Image Chain-of-Thought (MI-CoT)

Different versions of LLaVA were utilized to evaluate the performance of CoT-Prompting. For the model to extract information from all the images corresponding to a list of questions, we came up with a novel approach, namely a Multi-Image chain of thoughts (MI-CoT). Under this technique, the three images were stacked on top of each other. The rationale for employing multi-image prompting was driven by the anticipation that the Large Language Model (LLM) would effectively distinguish and identify the specific image to be utilized for each question within a single prompt. Consider the images corresponding to the two prompt questions X_p and X_q , and the image for the main question X_r . LLaVA utilizes the CLIP visual encoder to get the visual feature Z_v :

$$Z_v = g(X_v) \quad (1)$$

where

$$X_v = X_p \cdot X_q \cdot X_r \quad (2)$$

The filenames were passed as a list in the same order in which they were stacked. To make sure that the dimensions were correct for feeding the resultant concatenated image X_v into the CLIP encoder, the size of the images was reduced along one dimension

using an autoencoder after basic pre-processing (normalization and padding) of the images. A basic neural-network-based autoencoder was employed and was trained on the train split for this purpose.

5 Experiments

For evaluating the performance of the models, an 85/15 train-test split was used. We made sure that the percentage share of questions with options a, b, c, and d was roughly the same in both the training and testing datasets. This was especially important in the case of the training dataset to ensure no bias imposed by any option during the training process. We used accuracy as the primary metric for judging the performance of the models and rouge scores for evaluating the correctness of the reasoning.

5.1 Models

We conducted a variety of experiments with both text and multimodal LLMs to gauge the difference in performance that comes about due to the change in the modality. LLaMA2-7b and Mistral-7b are the current state-of-the-art open-source LLMs for textual input. These models were tested with text-only inputs. We use these LLMs to highlight the difference in the level of performance between fine-tuned models versus using them straight out of the box, aka through zero-shot prompting. For the ablation study, we also experimented with GPT-4, which is the current state-of-the-art model for multimodal question-answering.

LLaVA and LLaVA-1.5 being multimodal were provided with the figures along with the textual input. All the models were trained on A100 GPU and were fine-tuned for 5 epochs with a batch size of 8. Weighted Adam optimizer was utilised and the learning rate was set to $2e-4$.

We also experimented with different LoRA values in the case of the LLaVA-1.5 model. LoRA or Low-Rank Adaptation [30], is a method to represent the weight changes during the training process in lower-ranked matrices. This is especially useful while fine-tuning general-purpose LLMs, as it speeds up the training process. A lower LoRA rank means fewer parameters are learned during the adaptation process, however, it results in a faster training process as well. We tested the 7b (7 billion) and 13b (13 billion) variants of LLaVA which correspond to the number of learning parameters. The different LLaVA configurations also formed the basis of our comparison of the performance of (MI-CoT) Prompting. For fine-tuning, open-source base model checkpoints from huggingface were utilized.

6 Results and Discussion

6.1 Model Performance

The results of the experiments with their accuracy scores on the test dataset are listed in Table 2. Mistral-7b and LLaMA2-7b being text-only models only take into account the textual data which means that they are bound to miss critical information in some

Table 2: Performance of text-only and multimodal (MM) models. Model training specifications such as LoRA Rank and whether MI-CoT Prompting was used have been mentioned. All models were fine-tuned except for GPT-4, for which the answers were extracted using zero-shot prompting

| Model | MI-CoT | Modality | Accuracy | Rouge1 | Rouge2 | RougeL | LoRA Rank |
|---------------|--------|-----------|--------------|--------------|--------------|--------------|-----------|
| LLaMA2-7b | × | Text Only | 0.25 | 0.380 | 0.187 | 0.315 | 8 |
| Mistral-7b | × | Text Only | 0.428 | 0.460 | 0.256 | 0.391 | 8 |
| GPT-4 | × | MM | 0.331 | - | - | - | - |
| LLaVA-13b | × | MM | 0.293 | 0.551 | 0.383 | 0.501 | 64 |
| LLaVA-1.5 7b | × | MM | 0.533 | 0.712 | 0.579 | 0.676 | 64 |
| LLaVA-1.5 13b | × | MM | 0.527 | 0.672 | 0.532 | 0.634 | 64 |
| LLaVA-1.5 13b | × | MM | 0.531 | 0.621 | 0.490 | 0.586 | 128 |
| LLaVA-13b | ✓ | MM | 0.291 | 0.383 | 0.184 | 0.306 | 64 |
| LLaVA-1.5 7b | ✓ | MM | 0.354 | 0.496 | 0.343 | 0.444 | 64 |
| LLaVA-1.5 13b | ✓ | MM | 0.653 | 0.686 | 0.585 | 0.656 | 64 |
| LLaVA-1.5 13b | ✓ | MM | 0.716 | 0.677 | 0.582 | 0.650 | 128 |

questions. We observed an accuracy score of 25.95% and 42.83% for LLaMA2-7b and Mistral-7b, respectively. Thus, we conclude that text-only LLMs are not capable of providing the right answers for a large number of multimodal questions which require multiple steps with complex reasoning to reach the final answer.

LLaVA is a model that can potentially answer complex questions due to its ability to process images. While the older LLaVA version with 13 billion parameters exhibited a lower accuracy than Mistral-7b, LLaVA-1.5 was able to perform significantly better than Mistral-7b. The best performance was seen when LLaVA-1.5, trained with 13 billion parameters, was fine-tuned with a LoRA rank of 128 and employed Chain of Thought Prompting with an accuracy score of 71.65%. A higher LoRA rank means that the model can learn more parameters during fine-tuning which makes it ideal for task-specific situations, such as answering complex physics questions. LLaVA-1.5 13b performs better than the 7b variant with an equal LoRA rank of 64 when multi-image prompting was utilized. This is because the larger number of trainable parameters allowed the model to learn and generalize better.

Table 3: Performance of text-only LLMs using zero-shot prompting and fine-tuning

| Model | Task | Modality | Accuracy(in %) | Rouge 1 | Rouge 2 | Rouge L |
|------------|------------------------|-----------|----------------|--------------|--------------|--------------|
| LLaMA2-7b | Zero Shot Prompting | Text Only | 14.22 | 0.301 | 0.096 | 0.201 |
| | Supervised Fine-tuning | Text Only | 25.95 | 0.380 | 0.187 | 0.315 |
| Mistral-7b | Zero Shot Prompting | Text Only | 23.32 | 0.259 | 0.083 | 0.180 |
| | Supervised Fine-tuning | Text Only | 42.83 | 0.460 | 0.256 | 0.391 |

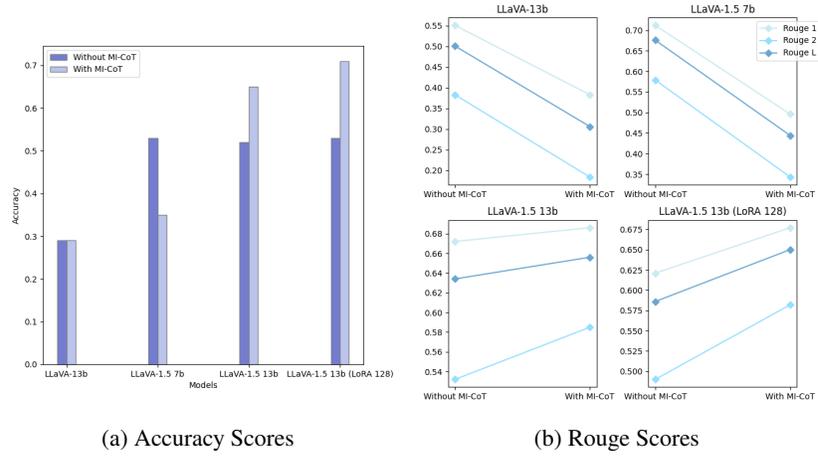


Fig. 4: Comparison of the accuracy and rouge scores of different LLaVA variants when trained using (MI-CoT) Prompting vs their non-CoT prompted supervised fine-tuned (SFT) counterparts

6.2 Zero Shot Prompting vs Supervised Fine-tuning

Table 3 shows the performance of LLaMA2-7b and Mistral-7b with zero-shot prompting and supervised fine-tuning. There is a marked improvement in the accuracy, Rouge1, Rouge2, and RougeL scores for both the models when fine-tuned on the dataset. This proves the assertion that current LLM models in their out-of-the-box configurations are not able to answer physics questions satisfactorily, and there is a need to fine-tune the models on domain-specific datasets to get better performance.

Zero-shot inferencing was done using the GPT-4 model. In most instances, GPT-4 failed to give correct answers and was not able to extract the entire information from the image. In some failure cases, GPT-4 needed more context than questions to make progress toward the solution.

6.3 Effect of Chain of Thought Prompting

For all variants of LLaVA-1.5 that were tested, there was an increase in the accuracy score when MI-CoT Prompting was employed as seen in Figure 4a except in the case of LLaVA-1.5 7b model. A smaller number of trainable parameters meant that the model was not able to process the more complex multi-image input, leading to a sharp dip in the performance. The difference was the most significant in the case of LLaVA-1.5 13b trained with LoRA as 128, which also gave the best performance out of all the models tested when trained using MI-CoT Prompting. The MI-CoT Prompting trained version also exhibited high rouge scores as seen in Table 2. It can be observed from Figure 4b that the rouge scores were higher in the LLaVA-1.5 13b CoT variants, showcasing the fact that models that were able to leverage the MI-CoT prompt also showed a bump in the reasoning capabilities. A marked improvement in all metrics, when multiple images were provided in the prompt in the case of LLaVA-1.5 13b variants, provides evidence

that the models were able to segregate and recognize the image that has to be used for each question present in a single prompt.

6.4 Error Analysis

Different types of errors were explored in [11]. We investigated the error cases that were thrown by the best-performing model, LLaVA-1.5 13b. Figure 5 shows the different types of errors that were encountered. Their descriptions are as follows:

- ◇ **Conceptual Error:** The model is not able to identify the concepts that have to be involved correctly. For instance, in Figure 5a, the model fails to identify that Kirchhoff’s loop rule has to be applied.
- ◇ **Grounding Error:** The model is able to identify the concept that has to be applied but does not formulate the correct equation. In the case of Figure 5b, the model is not able to apply the correct equation to get the centripetal acceleration of the cyclist.
- ◇ **Computational Error:** The model makes an algebraic mistake. In Figure 5c, the concept and the equations are correct, but the computation of the final answer is incorrect.

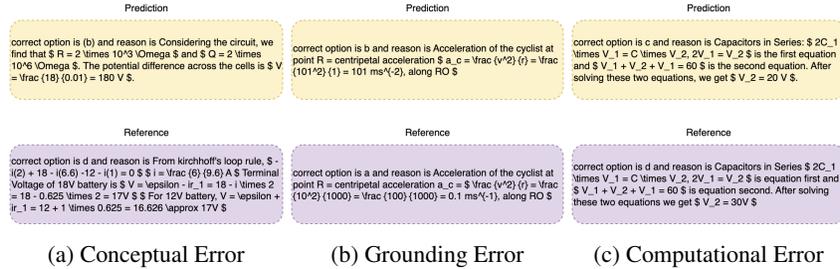


Fig. 5: Types of errors encountered by LLaVA-1.5 13b

7 Conclusion

This paper introduces the MM-PhyQA dataset, comprising high-quality problems solved by tested LLMs, serving as a benchmark for LLM performance in education. From our experiments, we concluded that the base configurations of Mistral-7b, LLaMA2, LLaVA-1.5, and GPT-4 struggled with complex reasoning tasks, but fine-tuning, particularly with MI-CoT prompting, showed promise, notably with the LLaVA-1.5 13b model. LLaVA’s image extraction abilities yielded high metric scores, and leveraging multimodality and MI-CoT Prompting, improved performance significantly. Future work may explore incorporating Reinforcement Learning from Human Feedback (RLHF) for model alignment and extending MI-CoT Prompting to other multimodal

tasks.

Acknowledgements. Rajiv Ratn Shah is partly supported by the Infosys Center for AI, the Center of Design and New Media, and the Center of Excellence in Healthcare at IIIT Delhi.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Touvron, Hugo, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov et al. "Llama 2: Open foundation and fine-tuned chat models." arXiv preprint arXiv:2307.09288 (2023).
2. Jiang, Albert Q., Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand et al. "Mistral 7B." arXiv preprint arXiv:2310.06825 (2023).
3. Anand, Avinash, Arnav Goel, Medha Hira, Snehal Buldeo, Jatin Kumar, Astha Verma, Rushali Gupta, and Rajiv Ratn Shah. "SciPhyRAG-Retrieval Augmentation to Improve LLMs on Physics Q &A." In International Conference on Big Data Analytics, pp. 50-63. Cham: Springer Nature Switzerland, 2023.
4. OpenAI. "GPT-4 Technical Report." ArXiv abs/2303.08774 (2023): n. pag.
5. Anand, Avinash, Kritarth Prasad, Ujjwal Goel, Mohit Gupta, Naman Lal, Astha Verma, and Rajiv Ratn Shah. "Context-Enhanced Language Models for Generating Multi-paper Citations." In International Conference on Big Data Analytics, pp. 80-94. Cham: Springer Nature Switzerland, 2023.
6. Anand, Avinash, Mohit Gupta, Kritarth Prasad, Ujjwal Goel, Naman Lal, Astha Verma, and Rajiv Ratn Shah. "KG-CTG: Citation Generation Through Knowledge Graph-Guided Large Language Models." In International Conference on Big Data Analytics, pp. 37-49. Cham: Springer Nature Switzerland, 2023.
7. Wei, Jason, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. "Chain-of-thought prompting elicits reasoning in large language models." *Advances in Neural Information Processing Systems* 35 (2022): 24824-24837
8. Lu, Pan, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. "Learn to explain: Multimodal reasoning via thought chains for science question answering." *Advances in Neural Information Processing Systems* 35 (2022): 2507-2521.
9. Zhang, Zhuosheng, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. "Multimodal chain-of-thought reasoning in language models." arXiv preprint arXiv:2302.00923 (2023).
10. Cobbe, Karl, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert et al. "Training verifiers to solve math word problems." arXiv preprint arXiv:2110.14168 (2021).
11. Arora, Daman, and Himanshu Gaurav Singh. "Have LLMs Advanced Enough? A Challenging Problem Solving Benchmark For Large Language Models." arXiv preprint arXiv:2305.15074 (2023). pag.
12. Welbl, Johannes, Nelson F. Liu, and Matt Gardner. "Crowdsourcing multiple choice science questions." arXiv preprint arXiv:1707.06209 (2017).

13. Wang, Xiaoxuan, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun and Wei Wang. "SciBench: Evaluating College-Level Scientific Problem-Solving Abilities of Large Language Models." ArXiv abs/2307.10635 (2023): n. pag.
14. Hendrycks, Dan, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. "Measuring massive multitask language understanding." arXiv preprint arXiv:2009.03300 (2020).
15. Huang, Yuzhen, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Jun-teng Liu et al. "C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models." arXiv preprint arXiv:2305.08322 (2023).
16. Chen, Jiaqi, et al. "GeoQA: A geometric question answering benchmark towards multimodal numerical reasoning." arXiv preprint arXiv:2105.14517 (2021).
17. Jin, Nengzheng, Joanna Siebert, Dongfang Li, and Qingcai Chen. "A survey on table question answering: recent advances." In China Conference on Knowledge Graph and Semantic Computing, pp. 174-186. Singapore: Springer Nature Singapore, 2022.
18. Masry, Ahmed, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. "ChartQA: A benchmark for question answering about charts with visual and logical reasoning." arXiv preprint arXiv:2203.10244 (2022).
19. Gupta, Deepak, Surabhi Kumari, Asif Ekbal, and Pushpak Bhattacharyya. "MMQA: A multi-domain multi-lingual question-answering framework for English and Hindi." In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018). 2018.
20. <https://openai.com/research/gpt-4v-system-card>
21. Driess, Danny, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid et al. "Palm-e: An embodied multimodal language model." arXiv preprint arXiv:2303.03378 (2023).
22. Liu, Haotian, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. "Visual instruction tuning." arXiv preprint arXiv:2304.08485 (2023).
23. Liu, Haotian, Chunyuan Li, Yuheng Li, and Yong Jae Lee. "Improved baselines with visual instruction tuning." arXiv preprint arXiv:2310.03744 (2023).
24. Radford, Alec, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry et al. "Learning transferable visual models from natural language supervision." In International conference on machine learning, pp. 8748-8763. PMLR, 2021.
25. Chen, Keqin, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. "Shikra: Unleashing Multimodal LLM's Referential Dialogue Magic." arXiv preprint arXiv:2306.15195 (2023).
26. Peng, Zhiliang, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. "Kosmos-2: Grounding Multimodal Large Language Models to the World." arXiv preprint arXiv:2306.14824 (2023).
27. Yao, Shunyu, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. "Tree of thoughts: Deliberate problem solving with large language models." arXiv preprint arXiv:2305.10601 (2023).
28. Besta, Maciej, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann et al. "Graph of thoughts: Solving elaborate problems with large language models." arXiv preprint arXiv:2308.09687 (2023).
29. <https://chat.openai.com/>
30. Edward J. Hu et al: LoRA: Low-Rank Adaptation of Large Language Models. CoRR abs/2106.09685 (2021)