

# “Improving Physics Reasoning in Large Language Models Using Mixture of Refinement Agents”

Raj Jaiswal<sup>\*1</sup>, Dhruv Jain<sup>\*2</sup>, Harsh Parimal Papat<sup>1</sup>, Abhishek Dharmadhikari<sup>1</sup>, Atharva Marathe<sup>1</sup>, Avinash Anand<sup>1</sup>, Shin’ichi Satoh<sup>3</sup>, Rajiv Ratn Shah<sup>1</sup>  
 MIDAS LAB, IIIT Delhi<sup>1</sup>, Indian Institute of Technology (BHU) Varanasi<sup>2</sup>, National Institute of Informatics, Tokyo<sup>3</sup>

## Motivation :

While analyzing physics problems and their CoT solutions generated with LLMs (Llama-3-70B & Gemma-2-27B), we observed three key errors made by them :

**Observation 1 :** LLMs in few cases struggle to fully grasp the objective of the question, along with misinterpreting the values of variables and constants provided in the question. Although this issue has been identified in only a few cases, it is significant one because it leads to solutions that fails to address the correct interpretation of a given question resulting in problem miscomprehension.

**Observation 2 :** LLMs struggle to apply the correct concepts or formulae with respect to the context of the given problem. This issue is a more recurring one in LLMs, especially for problems requiring considering a specific case rather than relying on a generic formula. For instance, the formula for calculating the moment of inertia varies depending on the distribution of mass.

**Observation 3 :** Many physics problems involve mathematical reasoning and algebraic computation, areas where LLMs tend to struggle. Computational errors account for the majority of errors in solutions generated by LLMs. LLMs struggles with accurate algebraic and arithmetic computations resulting in errors within the reasoning and final answer.

| Problem Miscomprehension  | Incorrect Concept Application   | Computational Error   |
|---|---|---|
| <b>Q:</b> From a uniform circular disc of radius $R$ and mass $5M$ , a small disc of radius $R/2$ is removed. The moment of inertia of the remaining . . . . .<br><b>R:</b> Let's break it down step by step:<br>1. The original disc has a uniform mass density, so its moment of inertia about an axis perpendicular to the plane and passing through the center:<br>$I = (1/2)MR^2$ , where $M$ is the total mass of the disc. | <b>Q:</b> The moment of inertia of a uniform cylinder of length $L$ and radius $R$ about its perpendicular bisector is 1. What is the ratio $L / R$ such that the moment of inertia is minimum?<br><b>R:</b> A nice physics question!<br>The moment of inertia of a uniform cylinder about its perpendicular bisector is given by:<br>$I = (1/2)MR^2$ , where $M$ is the mass of the cylinder.<br>Since the cylinder is uniform, its mass . . . . . | <b>Q:</b> A magnetic needle of magnetic moment $6.7 \times 10^{-2} \text{ Am}^2$ and moment of inertia $7.5 \times 10^{-6} \text{ kg m}^2$ is performing simple harmonic oscillations in a magnetic field of $0.01 \text{ T}$ . Time for 10 oscillations<br><b>R:</b> Let's break it down step by step:<br>. . . . .<br>Step 4: Calculate the time period<br>Substituting the values, we get, $T = 2\pi \sqrt{\frac{I}{MB}} = \sqrt{2} \times \sqrt{\frac{7.5 \times 10^{-6}}{6.7 \times 10^{-2} \times 0.01}} = 0.698 \text{ s}$<br>Step 5: Time for the 10 oscillations . . . . . |

## Dataset: PhysicsQA

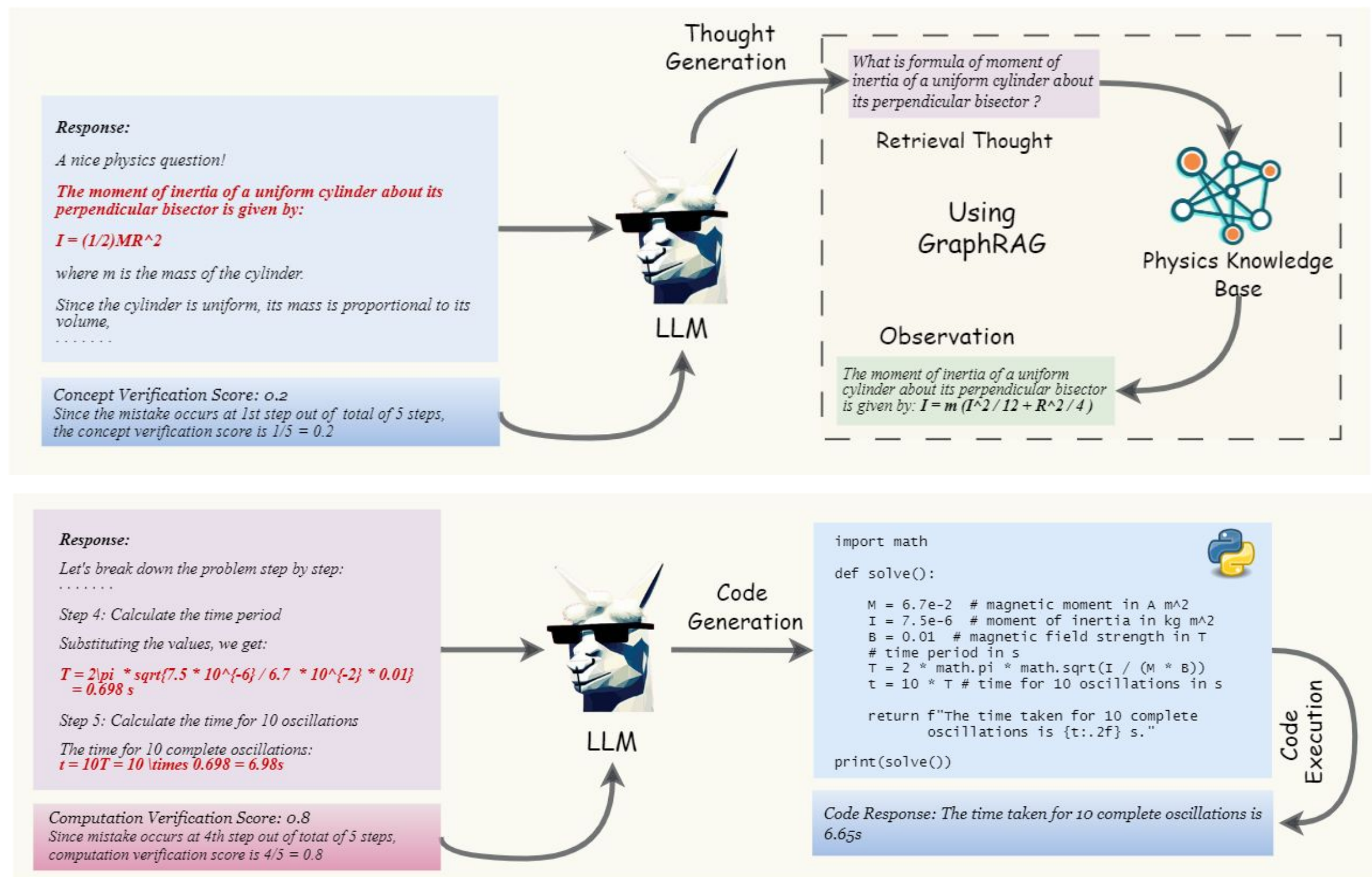
Benchmarks like MMLU, SciEval and ScienceQA focus on foundational knowledge and general reasoning, while more challenging ones like OlympiadBench and JEEBench require advanced reasoning skills. To bridge the gap, we curated our own dataset PhysicsQA, containing set of 370 diverse, intermediate level high school physics problems that provide a balanced challenge, allowing a exhaustive evaluation and step by step solution analysis of open-source LLMs on physics problems. Table 1 illustrates the topic-wise distribution of the questions, providing a clear overview of the areas covered.

| Chapter Name                       | Percentage |
|------------------------------------|------------|
| Electromagnetism                   | 29.8%      |
| Mechanics and Kinematics           | 21.8%      |
| Thermodynamics and Heat            | 15.7%      |
| Waves and Optics                   | 15.4%      |
| Nuclear and Modern Physics         | 8.9%       |
| Material Properties and Elasticity | 8.3%       |

Table 1: Topic-wise Distribution in PhysicsQA

## Mixture of Refinement Agents :

This introduces our mixture of refinement agents (MoRA) framework. We first discuss our motivation behind MoRA; then, we introduce the error identification stage and refinement agents. Finally, we discuss how these agents are routed iteratively to correct different errors in the solutions generated by the LLM.



### Algorithm 1: Error Identification and Iterative Refinement

**Require:** Question  $Q$ , Initial Solution  $S_0$ , GPT-4o  $\mathcal{L}$ , Refinement Agents  $\mathcal{R}$ , Maximum Iterations  $N$ , Threshold  $\epsilon$   
**Ensure:** Final refined solution to  $Q$   
 1:  $i = 0, S_i = S_0$   
 2: **while**  $i < N$  **do**  
 3:  $(F_{\text{obj}}^i, F_{\text{val}}^i, \text{Score}_{\text{concept}}^i, \text{Score}_{\text{comp}}^i) \leftarrow \mathcal{L}(Q, S_i)$   
 4: **if**  $F_{\text{obj}}^i == -1$  **or**  $F_{\text{val}}^i == -1$  **then**  
 5:  $S_{i+1} \leftarrow \mathcal{R}_{\text{miscomprehension}}(Q, S_i)$   
 6: **else if**  $\text{Score}_{\text{concept}}^i < 1 - \epsilon$  **then**  
 7:  $S_{i+1} \leftarrow \mathcal{R}_{\text{concept}}(Q, S_i)$   
 8: **else if**  $\text{Score}_{\text{comp}}^i < 1 - \epsilon$  **then**  
 9:  $S_{i+1} \leftarrow \mathcal{R}_{\text{computation}}(Q, S_i)$   
 10: **else**  
 11: **return**  $S_i$   
 12: **end if**  
 13:  $i \leftarrow i + 1$   
 14: **end while**  
 15: **return**  $S_N$

## Setup :

- Datasets :** In our experiments, we use four datasets: SciEval-Static, PhysicsQA, MMLU High School and MMLU College. SciEval-Static is a subset of SciEval, consisting 164 questions from physics divided into multiple sub-topics. MMLU, consists of a 118 College level and 173 high school multiple-choice questions from various disciplines.
- LLMs :** We utilize the API of a range of models with varying parameters and capabilities including LLaMa-3-70B, LLaMa 3.1-405B, Gemma-2-27B, Gemini-1.5-Flash, GPT- 3.5 Turbo and GPT-4 as our LLMs for the evaluation. We use same prompts for all the datasets and LLMs during our Evaluation.
- Baselines :** We employ an Answer-only approach (AO), where the model is given a question with four options and asked to select the correct answer without any explanation relying solely on its pre-existing knowledge. In contrast, few-shot prompting uses a few examples to help the model learn and apply that knowledge to similar tasks. Chain-of-Thought (CoT) prompting guides the model to generate intermediate reasoning steps, improving its performance on complex tasks by breaking them down into smaller, more manageable parts. These three approaches form our primary baselines.
- Evaluation :** Most of the existing works measure the mathematical reasoning quality of LLMs by directly comparing the final answer and calculating the overall accuracy on a given dataset. We choose to follow the same evaluation for physics reasoning as well.

| Model            | SciEval-Static |               |               | PhysicsQA     |               |               | MMLU - High |               |               | MMLU - College |               |               |
|------------------|----------------|---------------|---------------|---------------|---------------|---------------|-------------|---------------|---------------|----------------|---------------|---------------|
|                  | AO             | CoT           | 3-Shot        | AO            | CoT           | 3-Shot        | AO          | CoT           | 3-Shot        | AO             | CoT           | 3-Shot        |
| LLaMa-3-70B      | 70.07%         | 82.23%        | 63.41%        | 38.37%        | 56.76%        | 59.29%        | 60.16%      | 72.88%        | 73.66%        | 59.41%         | 71.76%        | 71.76%        |
| LLaMa 3.1 405B   | <b>79.87%</b>  | 89.63%        | <b>82.92%</b> | <b>50.81%</b> | 76.75%        | <b>78.37%</b> | <b>72%</b>  | 91.52%        | <b>88.98%</b> | <b>75.29%</b>  | <b>88.23%</b> | <b>85.29%</b> |
| Gemma-2-27B      | 60.36%         | 79.26%        | 53.04%        | 39.18%        | 54.59%        | 59.45%        | 55.93%      | 77.11%        | 74.45%        | 51.11%         | 73.52%        | 67.64%        |
| Gemini 1.5 Flash | 68.29%         | 85.97%        | 81.70%        | 44.86%        | 62.97%        | 69.72%        | 58.47%      | 79.66%        | 80.05%        | 60.58%         | 72.35%        | 72.94%        |
| GPT 3.5 Turbo    | 41.46%         | 66.46%        | 48.78%        | 28.10%        | 42.70%        | 42.70%        | 47.45%      | 58.47%        | 33.89%        | 35.29%         | 50.58%        | 42.35%        |
| GPT4o            | 64.02%         | <b>92.68%</b> | 81.09%        | 49.45%        | <b>79.45%</b> | <b>78.37%</b> | 62.71%      | <b>94.06%</b> | 87.28%        | 70%            | 84.70%        | 84.17%        |

Table 2: Experimentation of Answer-Only (AO), CoT and Few-Shot (3-shot) on different Datasets

| Model              | Dataset          | AO     | COT    | 3-Shot | MORA          |
|--------------------|------------------|--------|--------|--------|---------------|
| <b>Gemma 2 27B</b> | MMLU College     | 51.11% | 73.52% | 67.64% | <b>82.20%</b> |
|                    | MMLU High School | 55.93% | 77.11% | 74.45% | <b>75.88%</b> |
|                    | PhysicsQA        | 39.18% | 54.59% | 59.45% | <b>70.62%</b> |
|                    | SciEval-Static   | 60.36% | 79.26% | 53.04% | <b>88.76%</b> |
| <b>LLaMa 3 70B</b> | MMLU College     | 59.41% | 71.76% | 71.76% | <b>78.82%</b> |
|                    | MMLU High School | 60.16% | 72.88% | 73.66% | <b>78.81%</b> |
|                    | PhysicsQA        | 38.37% | 56.76% | 59.29% | <b>70.14%</b> |
|                    | SciEval-Static   | 70.07% | 82.23% | 63.41% | <b>86.58%</b> |

Table 3: Comparison of baseline approaches with MoRA across four datasets: SciEval-Static, PhysicsQA, MMLU High School and College based on final answer accuracy.

## Error Analysis :

- LLMs demonstrate good problem comprehension ability for physics question.
- Open source LLMs sometimes struggles to retrieve correct physics concept and formulae while reasoning.
- Open-source LLMs struggles with algebraic and arithmetic computation required while solving physics questions.

| Error Type                      | Dataset          | GPT-4o | Gemma 2-27B | LLaMa 3-70B |
|---------------------------------|------------------|--------|-------------|-------------|
| <b>Computational Error</b>      | MMLU College     | 2.54%  | 5.08%       | 9.32%       |
|                                 | MMLU High School | 2.35%  | 3.53%       | 6.47%       |
|                                 | PhysicsQA        | 8.92%  | 22.16%      | 21.08%      |
|                                 | SciEval-Static   | 3.06%  | 10.37%      | 10.98%      |
| <b>Problem Miscomprehension</b> | MMLU College     | 0.00%  | 0.85%       | 1.69%       |
|                                 | MMLU High School | 0.59%  | 0.59%       | 1.18%       |
|                                 | PhysicsQA        | 0.54%  | 2.16%       | 1.62%       |
|                                 | SciEval-Static   | 0.00%  | 1.22%       | 1.83%       |
| <b>Wrong Concept</b>            | MMLU College     | 0.85%  | 12.71%      | 10.17%      |
|                                 | MMLU High School | 3.53%  | 8.24%       | 11.76%      |
|                                 | PhysicsQA        | 7.57%  | 17.11%      | 18.92%      |
|                                 | SciEval-Static   | 1.22%  | 7.36%       | 9.15%       |

Table 4: Error Analysis of incorrect physics CoT solutions of different models across four datasets.

## Ablation :

- Problem miscomprehension errors are mitigated with simple instruction prompting and error feedback.
- Open-source LLM performers moderately in identifying the conceptual mistake and retrieval thought generation.
- Using code-driven refinement significantly corrects the computational errors.

| Error Type                         | Dataset          | Gemma 2-27B | LLaMa 3-70B |
|------------------------------------|------------------|-------------|-------------|
| <b>Computational Refinement</b>    | MMLU College     | 100%        | 81.8%       |
|                                    | MMLU High School | 33.3%       | 75.0%       |
|                                    | PhysicsQA        | 73.3%       | 72.6%       |
|                                    | SciEval-Static   | 57.1%       | 60.0%       |
| <b>Miscomprehension Refinement</b> | MMLU College     | 37.5%       | 33.3%       |
|                                    | MMLU High School | 16.7%       | 37.5%       |
|                                    | PhysicsQA        | 48.7%       | 46.9%       |
|                                    | SciEval-Static   | 62.5%       | 57.1%       |
| <b>Concept Refinement</b>          | MMLU College     | 100%        | 100%        |
|                                    | MMLU High School | 100%        | 100%        |
|                                    | PhysicsQA        | 62.5%       | 66.7%       |
|                                    | SciEval-Static   | 100%        | 66.7%       |

Table 5: Ablation studies for different refinement agent in MoRA using Gemma-2-27B and Llama-3-70B across four datasets, evaluated by refinement rate.