

# MODULAR REFINEMENT OF SMALL LANGUAGE MODELS FOR PHYSICS REASONING VIA LOCALIZED ERROR FEEDBACK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) excel at many reasoning tasks but struggle with scientific domains like physics, which demand precise mathematical calculations alongside deep conceptual and factual understanding. In complex physics problem-solving, LLMs commonly falter due to three core issues: misunderstanding the problem, incorrect application of concepts, and calculation mistakes. These challenges are more pronounced in small LLMs due to their limited capacity, making them more prone to failures. To address these limitations, we propose a modular reinforcement learning refinement framework tailored for small LLMs, integrating first-step error localization and correction through a Reinforcement Learning agent-guided feedback mechanism. We also introduce PhysicsQA, a diverse benchmark of 370 physics problems designed to evaluate LLM reasoning across the aforementioned dimensions. Using our Framework, experimental results demonstrate improvements up to 16% in final answer accuracy reasoning using Small language models over existing results.

## 1 INTRODUCTION

Scientific reasoning, particularly in the field of physics, requires a deep understanding that spans multiple disciplines. It demands not only domain-specific knowledge but also the integration of mathematical calculation with theoretical concepts, applying abstract principles and formulae across various contexts and scenarios. Unlike purely mathematical reasoning, which is abstract and symbolic, physics reasoning requires grounding those abstractions in real-world phenomena and physical laws. Successfully solving these challenges is a fundamental aspect of human intelligence, as it entails not just recalling information but adapting knowledge to solve diverse complex problems. Large language models (LLMs) are growing in size, but bigger is not always better Boye & Moell (2025). Solving complex reasoning problems still an open challenge for small open source LLMs Srivastava et al. (2025) models.

Chain of Thoughts (CoT) Wei et al. (2022b) can enable LLMs to solve complex reasoning tasks, it is highly sensitive to individual mistakes and vulnerable to error accumulation Shen et al. (2021). If a tiny mistake occurs, it can change the meaning of the whole statement Xiao et al. (2023), leading to incorrect answers Cobbe et al. (2021). Alternatively Retrieval-Augmented Generation (RAG) improves factual accuracy by incorporating external knowledge, but they primarily focus on information retrieval rather than reasoning correction. One approach to address these challenge can be collecting question and solution trajectory annotations and finetune LLMs to enhance these capabilities, similar to recent mathematical reasoning works Luo et al. (2023); Yuan et al. (2024). However, the process of such annotations and finetuning is time-consuming and costly. Finetuning may sometimes lead to a larger decline in CoT reasoning performance and can compromise the faithfulness of CoT reasoning Lobo et al. (2024). It also leads to challenges such as catastrophic forgetting McCloskey & Cohen (1989). Reinforcement Learning from Human Feedback (RLHF) is perhaps the most well-known application of RL techniques for finetuning LLMs. Reward models in RLHF often favor verbosity, leading models to generate unnecessarily long or redundant outputs Chiang & Lee (2024). However, despite these advances, RL-enhanced LLMs still rely primarily on internal knowledge and language modeling Wang et al. (2024). This reliance becomes a major limitation for time-sensitive or knowledge-intensive questions, where the model’s static knowledge

base may be outdated or incomplete, often resulting in inaccuracies or hallucinations. Agentic reasoning addresses these limitations by enabling LLMs to dynamically interact with both external resources and environments throughout the reasoning process Xiong et al. (2025); Patil & Jadon (2025). Open source small LLMs (SLMs) struggles to directly identify reasoning mistakes in their own solutions Li et al. (2024); Tyen et al. (2024), making them unreliable for self verification and refinement. SLMs lack the necessary depth Zhang et al. (2024) to reliably detect mistakes or rectify their reasoning processes independently.

Based on the limitations of above Techniques to improved and correct Scientific reasoning as mentioned, our key contributions are:

- **PhysicsQA Benchmark:** Introduced PhysicsQA, a dataset of 370 diverse, intermediate-level high school physics problems with verified chain-of-thought solutions to evaluate LLM reasoning.
- **Error Taxonomy:** Identified three core LLM error types in physics problem solving: Problem Miscomprehension, Incorrect Concept Application, and Calculation Errors.
- **Agent Guided Feedback Training:** Designed a modular training pipeline using specialized agents to provide localized, error-specific feedback based on the identified error types.
- **Reinforcement Learning Framework:** Built a reinforcement learning framework for small LLMs that leverages error localization and structured feedback to improve physics reasoning.
- **Empirical Gains:** Achieved substantial accuracy improvements on multiple benchmarks, outperforming existing methods for small open-source LLMs.

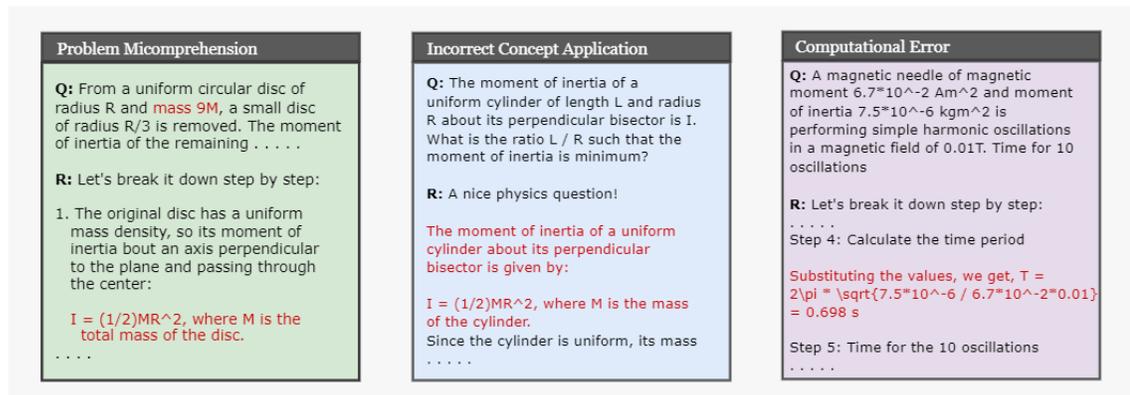


Figure 1: The illustration of three key error observations in the CoT solution of SLMs for physics problems (highlighted in red). (a) showcases problem miscomprehension, where the SLM response uses the incorrect value of variables given in the question here,  $M$  instead of  $9M$ , (b) showcases incorrect concept application in the SLM response, here incorrect moment of inertia formula for uniform cylinder, (c) demonstrate Calculation error within SLM response here, incorrect calculation of time period.

## 2 RELATED WORKS

### 2.0.1 LLM REASONING FOR PHYSICS

Researchers have begun exploring the potential of LLMs as reasoning tools in the physics domain Anand et al. (2024); Ding et al. (2023); Pang et al. (2024). Studies have demonstrated that LLMs can solve complex word problems requiring calculation and inference, often achieving near human-level accuracy, especially with effective prompting techniques such as few-shot learning using similar examples Ding et al. (2023), leveraging RLHF Anand et al. (2024) or implementing agentic system Pang et al. (2024). While much of this research focuses on general physics reasoning. LLMs have

108 been successfully applied to address multi-step reasoning tasks by generating intermediate reasoning  
109 steps, referred to as CoT, Auto-CoT Zhang et al. (2022), and Complex-CoT Fu et al. (2022), among  
110 others. LLMs tends to struggle with arithmetic calculations when solving math problems Gao et al.  
111 (2023), Leveraging the strengths of GPT4 Code Interpreter Zhou et al. (2023) has been integral to  
112 frameworks like MathCoder Wang et al. (2023), which is designed to improve the mathematical  
113 calculations capabilities of open-source models.

#### 114 2.0.2 LLM REASONING WITH FINE TUNING AND RLHF

115 While large LLMs exhibit excellent performance on mathematical reasoning tasks, adapting smaller  
116 models for these tasks remains an open problem Wei et al. (2022a). Some methods used ques-  
117 tions from existing training datasets and used prompting to generate solutions for fine tuning (FT)  
118 smaller models Ho et al. (2022); Fu et al. (2023). Others use various techniques to rephrase the  
119 questions to create more examples Vu et al. (2020) or multiple views of solutions Liang et al. (2023)  
120 to achieve better reasoning performance. Reinforcement Learning from Human Feedback (RLHF)  
121 Ouyang et al. (2022); Glaese et al. (2022) most often works by training a reward model to capture  
122 human preferences over a task. Full Step DPO Xu et al. (2025), a novel framework for mathematical  
123 reasoning that optimizes each step in the entire reasoning chain using step wise rewards.  
124

#### 125 2.0.3 LLM REASONING WITH EXTERNAL DATABASE

126 Lewis et al. (2020) proposed Retrieval-Augmented Generation (RAG) framework, which incorpo-  
127 rates a retrieval component to fetch relevant information from a given knowledge base. Integrating  
128 LLMs with knowledge representation tools, such as knowledge graphs (KGs) Mruthyunjaya et al.  
129 (2023), has further enhanced reasoning capabilities. Yao et al. (2024) demonstrated that augmenting  
130 LLMs with comprehensive external knowledge from KGs can significantly improve their perfor-  
131 mance and facilitate more robust reasoning processes.  
132

#### 133 2.0.4 SELF VERIFICATION WITH LLMs

134 Recent works Cobbe et al. (2021); Ling et al. (2024) have attempted to address the challenge of  
135 error detection in step-by-step reasoning. Miao et al. (2023) proposes using the LLM itself to verify  
136 the conditional correctness of each step in the reasoning chain, similar to how a human reviews  
137 their work. Accurate error recognition and correction are crucial for enhancing problem solving  
138 capabilities, as demonstrated by Li et al. (2024), which defines tasks to assess LLMs’ mathematical  
139 reasoning abilities in error identification and correction.  
140

### 141 3 DATASET : PHYSICSQA

142 Benchmarks like MMLU Hendrycks et al. (2020), SciEval Sun et al. (2024), focus on foundational  
143 knowledge, while more challenging ones like OlympiadBench He et al. (2024) and JEEBench Arora  
144 et al. (2023) require advanced reasoning skills. To bridge the gap, we curated our own dataset  
145 PhysicsQA, containing the set of diverse, intermediate-level high school physics problems that pro-  
146 vide a balanced challenge, allowing a exhaustive evaluation and analysis of LLMs on physics prob-  
147 lems. It has 370 high-school Indian Engineering Exam (JEE) physics problems sourced from PW  
148 Live (2024); Doubtnut that offer diversity in both topic coverage and difficulty. These problems are  
149 notably challenging, often requiring the application of multiple concepts, intricate calculations, and  
150 multihop reasoning. Each problem is accompanied by a correct CoT solution, verified using GPT4  
151 with human annotators in the loop. This enables researchers to assess physics reasoning beyond final  
152 answer accuracy and conduct rigorous, step by step evaluations of LLMs conceptual understanding  
153 and calculation abilities. Constructions details and chapters breakdown illustrated in 3  
154  
155

### 156 4 METHODOLOGY

157 In this section, we present our Agent-Guided Reinforcement Learning framework for enhancing  
158 physics reasoning for SLMs. Our approach begins with a taxonomy of fundamental error types  
159 commonly observed in complex physics problem solving. Building on this foundation, As illustrated  
160 in Figure 2, we introduce a structured RL training pipeline composed of three core components: (1)  
161

an error localization model that identifies the earliest mistake in the model’s reasoning; (2) An agentic feedback generator that produces localized, error-specific guidance; and (3) a reinforcement learning loop that fine-tunes the model to revise its reasoning using the provided feedback. In the following subsections, we detail each of these components and how they interact to enable effective physics reasoning in SLMs.

#### 4.1 TRAINING DATASET

We curated a finetuning dataset comprising 2,494 high school-level physics questions sourced from standard materials Pandey; Tipler (1999); Pinsky (1989); Halliday et al. which is different source from our benchmark PhysicsQA. Each question is paired with a detailed CoT solution. The dataset provides a diverse and challenging set of physics problems, accompanied by oracle level step by step reasoning, following the same process used for PhysicsQA. Training Sample shown in fig.19

#### 4.2 LORA SLM TRAINING WITH WARMUP

To make the RL training effective, we first warm-start the SLM through supervised finetuning, enabling it to internalize core reasoning structures and domain-specific solution patterns as shown in fig 2. For the subsequent reinforcement learning stage, we adopt LoRA adapters, which substantially reduce memory and compute overhead while training with our objective. This combination of warm initialization and parameter-efficient finetuning enables our framework to teach strong physics reasoning capabilities to small LLMs without full model finetuning.

#### 4.3 PHYSICS ERROR TAXONOMY

Based on a systematic analysis of solutions generated by wide range LLMs for complex physics problems, we identify three recurring categories of physics reasoning errors, as shown in Fig.1.

**Problem Miscomprehension**, SLMs 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. These misinterpretations result in fundamentally flawed reasoning trajectories that do not address the intended problem. **Conceptual Misapplication**, SLMs 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 SLMs, especially for problems requiring considering a specific case rather than relying on a generic formula. **Calculation Errors**, Many physics problems involve mathematical reasoning and algebraic calculation, areas where SLMs tend to struggle. They often make mistakes in these calculations, which propagate through the solution and affect both intermediate reasoning and final answers.

#### 4.4 ERROR LOCALIZATION

For error localization, we utilize LLaMa 3.1 405B. Prior work by Zhang et al. (2024) has shown that GPT4o can effectively use gold labels from initial solutions as signals for self-refinement. This oracle verifier setup serves as an upper bound on the performance of refinement agents. Inspired from this framework, we prompt LLaMa 3.1 405B to identify and localize errors in the solution, along with reasoning about the cause of each error. Our prompt based verifier performs three tasks: (a) Identify the first error step in the model-generated solution where model deviates from the correct reasoning. (b) Classify the error into one of three categories: (c) Problem Miscomprehension, Conceptual Error, or Calculation Error. Explain the error briefly, describing what went wrong in the first error reasoning step.

We utilize LLaMa 3.1 405B solely for error identification as shown Figure 2, which is used in our reward modeling and informs the routing of the error step to the appropriate feedback agent. It helps each agent to understand the initial point of failure, enabling effective and focused feedback for that stage.

#### 4.5 AGENTIC FEEDBACK GENERATION

We developed our feedback generator as shown in Figure 2, which produces structured, step-level refinement feedback to guide SLMs in correcting localized reasoning errors. It consists of a set of

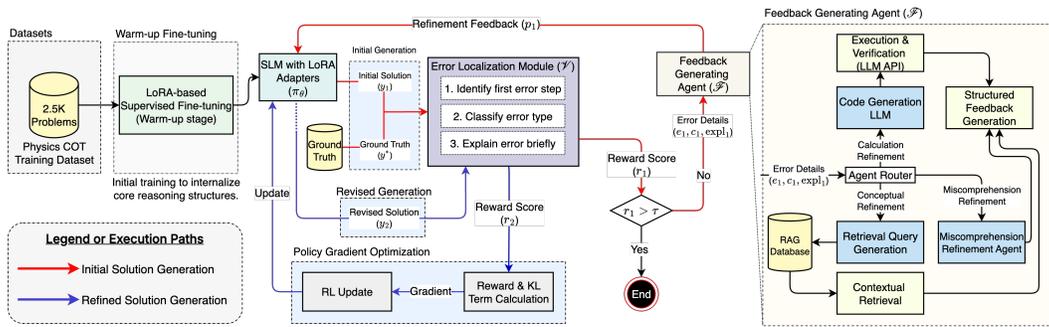


Figure 2: The architecture of our Agent-Guided Reinforcement Learning (AG-RL) framework. The process begins with a supervised fine-tuning warm-up stage. Subsequently, the SLM enters the iterative RL loop, starting with the Initial Policy Execution (red arrow) to generate a preliminary solution ( $y_1$ ). This solution is evaluated by the Feedback Generating Agent (F), which localizes the first error and generates structured feedback. This feedback informs a policy gradient update, triggering the Updated Policy Execution (blue arrow) to produce a revised solution ( $y_2$ ). A reward score, calculated from this revised solution, drives the final policy update.

specialized refinement agents, each designed to address one of the three core challenges in physics reasoning. The generated feedback is then incorporated into the model’s reasoning process to revise and improve its original solution attempt.

Miscomprehension identified in SLM generated solutions, can be corrected through instruction prompting. The agent handles such cases by incorporating the provided error explanation into a targeted feedback prompt, instructing the model to revise its understanding of the problem and generate a fully corrected solution in the next attempt.

To address incorrect concepts and misapplied formulae in the identified error step, we leverage RAG to produce targeted, context-aware corrections. Given a physics question, an erroneous reasoning step, and an explanation of the conceptual mistake, the agent performs the following sequence:

**Retrieval Query Generation**, The agent generates a focused query using the question, error step, and explanation. This query is optimized to retrieve the most relevant conceptual content from a domain-specific corpus. **Contextual Retrieval**, Using a vector store built with sentence-transformer embeddings, the agent retrieves key conceptual content required to correct the identified error step. **Structured Feedback Generation**, The retrieved context and original error description are used to generate structured, corrective feedback. This feedback is designed to help the model understand its mistake and apply the necessary correction in its reasoning.

We adopt a code generation approach to refine calculation and mathematical errors in the identified erroneous step. By leveraging executable code as a verification mechanism, the agent ensures accurate calculation correction and interpretability in the reasoning. Given a physics question, an erroneous step, and an explanation of the mistake, that guides small LLMs toward step-level correction. The agent follows a three-step process: **Code Generation**, The agent generates Python code that correctly performs the intended calculation, including proper variable initialization, unit conversions, and mathematical logic. **Execution & Verification**, The generated code is executed to validate correctness and produce a definitive numerical result. **Feedback Generation**, Based on the code and its output, the agent generates natural language feedback that explains the correct logic, clarifies the original error, and ensures unit consistency.

#### 4.6 AGENT GUIDED RL TRAINING

Self Correction remains a highly desirable yet largely ineffective capability in current LLMs Kamoi et al. (2024). To address this Kumar et al. (2024) introduced SCoRe, a multi-turn reinforcement learning (RL) framework aimed at teaching models to revise their own mistakes. This enables the model to explore diverse reasoning trajectories and potentially converge on more robust solutions.

270 However, RL based training pipelines are often complex and computationally demanding. Sidahmed  
 271 et al. (2024); Wang et al. (2025) demonstrate that integrating LoRA Hu et al. (2022) into RL setups  
 272 allows for efficient adaptation of reasoning patterns while retaining core model knowledge signifi-  
 273 cantly reducing training time, memory usage, and overall cost.

274 Building on these foundations, we introduce an efficient agent-guided RL training pipeline as shown  
 275 in figure 2, that improves physics reasoning in SLMs by refining their CoT reasoning using localized  
 276 error and agentic feedbacks.

277 Given a physics problem  $x$ , the model generate a first attempt solution  $y_1 \sim \pi_\theta(\cdot | x)$ . This  
 278 solution is passed to the error localization model, which identifies the first erroneous step in the  
 279 model reasoning, categorizes the type of error (miscomprehension, conceptual, or calculation), and  
 280 provides a brief explanation. If the model’s first response is deemed sufficiently accurate quantified  
 281 by a reward score  $r_1 > 0.95$ , the attempt is considered correct and skipped from the correction  
 282 loop. Otherwise, the feedback-generating agent uses the localized error information to construct a  
 283 correction feedback prompt  $p_1$ , which is given as the input to the model to produce a revised attempt  
 284  $y_2 \sim \pi_\theta(\cdot | x, y_1, p_1)$ . The agentic feedback assist the SLM to explore correct and robust reasoning  
 285 trajectory in the revised attempt  $y_2$ .

#### 287 4.7 REWARD DESIGN

288 The reward function is designed to encourage SLM to progressively refine the first error step in its  
 289 reasoning and is defined as:

$$291 \quad r = \frac{e_{\text{first}}}{n + 1}$$

292 where  $e_{\text{first}}$  is the first error step number and  $n$  is the total number of steps in the generated solution.  
 293 The first error step here is identified using the error localization module. This reward signal penalizes  
 294 earlier mistakes and encourages the model to push the first error further along the reasoning chain.  
 295 As training progresses, the model learns to push and ultimately eliminate its first error, thereby  
 296 producing increasingly accurate step-by-step solutions.

297 Given the revised attempt  $y_2$ , we evaluate it using our reward function resulting in the reward score  
 298  $r_2$ . The training objective is to shift the first error further down the reasoning chain or completely  
 299 eliminate it in  $y_2$  using agentic feedback and is given by:

$$302 \quad \max_{\theta} \mathbb{E}_{x, y_1, y_2} [r_2(y_2, y^*) - \beta_2 D_{\text{KL}}(\pi_\theta(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

303 Here,  $r_2$  reflects how much the model response has improved after agentic feedback. A higher  $r_2$   
 304 indicates a successful refinement, such as correcting an earlier mistake or pushing the first error  
 305 further in the reasoning chain. The KL regularization term ensures the updated policy remains close  
 306 to a reference policy during training, promoting stable optimization. This setup allows the model to  
 307 iteratively improve its reasoning capabilities over training steps, guided by localized, high-quality  
 308 agentic feedback.

#### 311 4.8 TRAINING SETUP

312 The warmup stage finetuning was done for 2 epoch using a single H100 GPU using AdamW opti-  
 313 mizer and a learning rate of  $5e-5$  and a batch size of 4. For agentic RL training, we utilized  
 314 Distributed Data Parallel in PyTorch, using AdamW optimizer, learning rate of  $5e-6$ , and a batch  
 315 size of 4 per GPU. The RL training was performed for 6 epochs using 4 H100 GPUs.

## 318 5 EXPERIMENTS

### 320 5.0.1 BENCHMARK DATASETS

321 In our experiments, we used four datasets: SciEval-Static, MMLU High School and College,  
 322 JEEBench and PhysicsQA. SciEval-Static is a subset of SciEval , consisting 164 questions from  
 323 physics divided into multiple sub-topics. MMLU, consists of a 118 College level and 173 high

LLM Models	SciEval-Static			MMLU-High			MMLU-College			JEEBench			PhysicsQA		
	AO	CoT	3-Shot	AO	CoT	3-Shot	AO	CoT	3-Shot	AO	CoT	3-Shot	AO	CoT	3-Shot
<b>Ultra Small</b>															
Qwen 2.5 1.5B Instruct	50.60	62.73	53.65	32.94	47.05	31.12	45.45	52.62	42.23	46.34	37.53	41.46	27.29	30.64	34.59
LLaMa 3.2 1B Instruct	40.24	44.51	35.97	30.00	30.58	32.34	26.36	33.63	29.32	30.89	32.52	34.95	20.27	21.35	22.16
LLaMa 3.2 3B Instruct	45.12	62.26	48.17	28.23	50	44.32	32.72	53.63	43.65	43.90	30.21	40.65	27.02	27.67	40.81
Phi 3.5 mini 3.8B Instruct	54.87	62.43	62.19	48.23	55.17	54.67	41.81	65.63	60.50	43.08	35.90	43.08	32.70	33.35	33.51
<b>Small</b>															
LLaMa 3.1 8B Instruct	48.17	79.26	45.73	43.52	57.64	51.26	40.90	61.81	63.35	39.83	38.21	39.83	27.02	35.67	33.24
OLMo 7B Instruct-hf	42.68	39.63	44.51	26.47	28.82	29.83	36.36	28.18	35.26	34.14	32.52	31.23	12.97	22.16	21.89
Phi 3 medium 14B Instruct	59.14	79.87	65.85	54.11	67.64	58.75	52.72	80.90	68.37	47.15	39.02	39.02	33.51	51.35	54.32
<b>Intermediate</b>															
Qwen2.5 32B Instruct	78.04	90.85	81.70	68.82	89.41	83.25	63.63	91.81	83.25	51.65	57.72	60.16	53.24	75.20	78.08
Qwen 2.5 72B Instruct	71.95	98.29	76.82	71.76	88.23	75.63	67.27	90.00	84.56	50.40	62.60	58.53	54.59	75.94	81.08
LLaMa 3 70B Instruct	70.73	91.46	87.80	62.35	85.29	78.35	60.90	90.90	87.89	53.65	58.53	65.85	43.78	76.21	74.05
Gemma 2 27b Instruct	56.09	82.92	61.58	53.52	70.00	69.18	56.36	77.27	78.29	45.52	47.96	44.71	37.02	52.97	63.78
<b>Large</b>															
LLaMa 3.1 405B Instruct	80.48	87.80	84.14	73.52	85.29	79.36	78.00	87.27	85.43	62.50	68.33	70.00	51.08	74.32	70.00
Mixtral 8x7B Instruct v0.1	53.25	64.37	65.62	46.47	61.17	62.87	49.09	53.63	52.31	29.16	30.00	37.50	34.32	38.91	43.78
Mixtral 8x22B Instruct v0.1	59.37	64.37	65.62	51.76	72.94	73.65	46.36	80.00	65.43	45.83	47.50	48.33	37.02	55.94	59.72
Deepseek R1 685B	53.25	61.28	65.38	50.12	58.35	60.12	45.34	55.25	56.12	36.13	44.35	51.65	39.13	44.28	43.15
<b>Proprietary</b>															
GPT 4o	64.02	92.68	81.09	62.71	94.06	87.28	70.00	84.70	84.17	51.67	66.67	60.83	49.45	79.45	78.37
Gemini 1.5 Flash	68.29	85.97	81.70	58.47	79.66	80.05	60.58	72.35	72.94	40.00	61.66	59.87	44.86	62.97	69.72

Table 1: We report final answer accuracy (%) of different prompting strategies Answer-Only (AO), Few-shot (3 shot), and Chain of Thought (CoT) across five Evaluation Physics benchmark Dataset. Models span a wide range of scales and architectures, including instruction-tuned LLM model variants. Accuracy highlighted in green indicates best performance.

Ultra Small	SciEval-Static					MMLU-High					MMLU-College				
	CoT	RAG	FT	DPO	Ours	CoT	RAG	FT	DPO	Ours	CoT	RAG	FT	DPO	Ours
Qwen 2.5 1.5B Instruct	62.73	68.12	57.93	59.32	79.29	47.05	59.41	41.18	46.47	55.03	52.62	62.27	43.22	43.22	53.64
LLaMa 3.2 1B Instruct	44.51	57.37	55.49	56.93	68.09	30.58	38.02	42.35	39.41	55.29	33.63	45.82	38.14	36.44	57.89
LLaMa 3.2 3B Instruct	62.26	66.87	61.59	57.93	81.52	50	55.29	48.29	47.65	67.20	53.63	61.82	45.76	50.85	69.09
Phi 3.5 mini 3.8B Instruct	62.43	61.01	59.17	56.87	65.19	55.17	60.53	55.97	56.12	65.83	65.63	67.27	62.77	64.34	73.12

Table 2: Final answer accuracy (%) across Physics Benchmarks for foundational knowledge using different strategies : Vanilla Chain-of-Thought prompting (CoT), Retrieval-Augmented Generation (RAG), supervised finetuning (FT), Direct Preference Optimization (DPO), and proposed framework (Ours). Accuracy highlighted in green indicates best performance.

school multiple-choice questions from various disciplines. JEEBench consists of 123 questions from Physics. Our PhysicsQA comprises 370 carefully selected complex high school physics questions sourced from online resources.

### 5.0.2 LLMs

We conduct experiments across a diverse set of open and closed-source LLMs, ranging from 1.5B to 405B parameters, to evaluate their physics reasoning capabilities on benchmark. The models, including Qwen2.5 1.5B, 32B and 72B Team (2024), LLaMA 3.1 8B , 70B, 405B, LLaMA 3.2 1B, 3B and LLaMA 3.3 70B Grattafiori et al. (2024), Phi 3.5 3.8B mini and Phi 3 14B medium Microsoft

Ultra Small	JEEBench					PhysicsQA				
	CoT	RAG	FT	DPO	Ours	CoT	RAG	FT	DPO	Ours
Qwen 2.5 1.5B Instruct	37.53	40	39.02	30.08	54.86	30.64	36.91	23.51	24.32	49.62
LLaMa 3.2 1B Instruct	32.52	35.50	28.46	35.77	49.17	21.35	28.64	23.24	18.92	39.02
LLaMa 3.2 3B Instruct	30.21	40.83	36.59	38.21	57.34	27.67	31.21	27.84	25.95	46.73
Phi 3.5 mini 3.8B Instruct	35.90	41.25	38.04	41.65	50.23	33.35	41.59	39.19	41.49	53.22

Table 3: Final answer accuracy (%) across Complex Physics Reasoning Benchmarks using different strategies : Vanilla Chain-of-Thought prompting (CoT), Retrieval-Augmented Generation (RAG), supervised finetuning (FT), Direct Preference Optimization (DPO), and proposed framework (Ours). Accuracy highlighted in green indicates best performance.

Research (2024), OLMo-7B Groeneveld et al. (2024), Gemma 2-27B Team et al. (2024b), Mixtral-8x7B and Mixtral-8x22B Jiang et al. (2024). We also include GPT4o Hurst et al. (2024) and Gemini 1.5 Flash Team et al. (2024a) as representative closed-source models.

## 5.1 SETUP

### 5.1.1 PROMPTING STRATEGIES

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 Yasunaga et al. (2023) uses a few examples to help the model learn and apply that knowledge to similar tasks. 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.

### 5.1.2 RETRIEVAL AUGMENTED GENERATION (RAG)

Unlike vanilla models that rely solely on intrinsic knowledge, our RAG framework using Langchain enhances reasoning by incorporating domain specific context by external physics knowledge PDFs contains formulae only required to solve ?. PDFs are chunked into 500 character segments with a 50-character overlap to preserve context and embedded using OpenAI’s `text-embedding-ada-002`. During inference, the top-k relevant chunks are retrieved and integrated into a structured prompt instructing the model to “Think step-by-step” using the provided context. This approach explicitly injects relevant formulas and concepts and improving solution accuracy in physics reasoning.

### 5.1.3 SLM FINETUNING

We finetuned ultra small models on our finetuning dataset consists of 2,494 high school-level physics questions sourced from standard Indian JEE preparation materials ?, each paired with a detailed CoT solution. The models include Qwen 2.5 1.5B Instruct, LLaMa 3.2 1B and 3B Instruct, Phi 3.5 mini 3.8B Instruct. finetuning followed a CoT approach, where each input instruction was framed as: *”You are an expert physics assistant. You are given a question. Your task is to generate the final solution of the given question. Let’s think step by step.”*. We utilized Hugging Face’s SFT library with the AdamW optimizer, a learning rate of  $5e-5$ , and a batch size of 2. All models were trained for 3 epochs on a single H100 GPU, with training time ranging from 45 to 60 minutes per model and employed LoRA adapters via the PEFT framework.

### 5.1.4 DIRECT PREFERENCE OPTIMIZATION (DPO)

Starting from our SFT-trained base small LLM model, the pairs consisting of a preferred (correct) and rejected (incorrect) response were used to train the model using the DPO loss, which encourages higher likelihood for preferred responses without requiring on-policy sampling during training. We employed Hugging Face’s `trl DPOTrainer` for this stage. Training was performed on a single H100 GPU using a batch size of 2, gradient accumulation steps of 4, and a learning rate of  $5e-5$  for 3

epochs. Each training run took approximately 1.5 to 2 hours. The optimizer used was AdamW, and early stopping was employed based on validation loss.

### 5.1.5 EVALUATION

Luo et al. (2023) 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. All models are evaluated using final answer accuracy. we further perform a step-by-step comparison when the final answer is incorrect. Using ground truth reasoning traces as reference Xia et al. (2025), we analyze each intermediate step in the model’s CoT output to identify specific reasoning error types.

## 6 RESULTS

In Table 1, Ultra small models perform poorly across all benchmarks and prompting strategies, reflecting limited foundational knowledge and a high tendency to hallucinate across all benchmarks. This underscores the complexity of CoT reasoning tasks where models still struggle.

In Table 2 and 3, we evaluate Ultra small model using vanilla CoT reasoning across various prompting strategies. RAG based prompting consistently improved accuracy by providing relevant external context. Surprisingly, FT degraded performance. Analysis suggests that models mimicked the surface form of complex CoT examples without internalizing the underlying logic, leading to flawed intermediate steps and incorrect answers. The fixed solution trajectories in the FT data further limited generalization. Supervised FT alone proved insufficient for teaching robust domain reasoning, especially in smaller models. DPO also failed to improve results, likely because it relies on a competent base model, which in our case lacked reliable intermediate reasoning.

Our approach consistently improves accuracy highlighted in table 2 and 3, across diverse benchmarks and all baselines. Our refinement agents are able to correct early stage errors. Notably, improvements in Ultra small models, suggesting that lightweight, reward guided corrections using agentic feedback can offset capacity limitations. However, in a few cases, performance drops marginally due to over corrections or agent misrouting, especially when initial answers are partially correct but the verifier triggers unnecessary edits. These fluctuations highlight the sensitivity of step level reward modeling and the importance of accurate error localization.

## 7 CONCLUSION

We tackle the underexplored problem of physics reasoning in SLMs, which often fail due to miscomprehension, conceptual, and calculation errors. Existing methods like FT and DPO rely on broad data level alignment but struggle with fine grained, step level corrections. Our core contribution is a lightweight, modular framework that uses agentic feedback and reinforcement learning, guided by a first error step reward. Unlike FT and DPO, our method avoids costly retraining and enables precise, iterative refinement through specialized agents.

## 8 LIMITATION AND FUTURE WORK

One major challenge in building our error localization module was designing a custom verifier to assess SLM generated physics reasoning. Prior studies Ma et al. (2025) emphasize that such verifiers demand large, high quality CoT solutions to perform reliably. Currently, we use LLaMa 3.1 405B as an oracle verifier. Future work includes a lightweight, domain-adapted verifier using a finetuned SLM.

## REFERENCES

Avinash Anand, Kritarth Prasad, Chhavi Kirtani, Ashwin R Nair, Mohit Gupta, Saloni Garg, Anurag Gautam, Snehal Buldeo, and Rajiv Ratn Shah. Enhancing llms for physics problem-solving using reinforcement learning with human-ai feedback, 2024. URL <https://arxiv.org/abs/2412.06827>.

- 486 Daman Arora, Himanshu Gaurav Singh, et al. Have llms advanced enough? a challenging problem  
487 solving benchmark for large language models. *arXiv preprint arXiv:2305.15074*, 2023.  
488
- 489 Johan Boye and Birger Moell. Large language models and mathematical reasoning failures. *arXiv*  
490 *preprint arXiv:2502.11574*, 2025.
- 491 Cheng-Han Chiang and Hung-yi Lee. Over-reasoning and redundant calculation of large language  
492 models. *arXiv preprint arXiv:2401.11467*, 2024.  
493
- 494 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
495 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
496 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 497 Jingzhe Ding, Yan Cen, and Xinyuan Wei. Using large language model to solve and explain physics  
498 word problems approaching human level, 2023. URL [https://arxiv.org/abs/2309.](https://arxiv.org/abs/2309.08182)  
499 08182.
- 500 Doubtnut. DC Pandey Physics - Questions, Answers, Solutions. [https://www.doubtnut.](https://www.doubtnut.com/books/class-11-dc-pandey-physics-english-medium-in-hindi-download-questions-and-answers)  
501 [com/books/class-11-dc-pandey-physics-english-medium-in-hindi-download-questions-and-answers](https://www.doubtnut.com/books/class-11-dc-pandey-physics-english-medium-in-hindi-download-questions-and-answers)  
502 Accessed: 2025-07-31.  
503
- 504 Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting  
505 for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*,  
506 2022.
- 507 Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. Specializing smaller language  
508 models towards multi-step reasoning. In *International Conference on Machine Learning*, pp.  
509 10421–10430. PMLR, 2023.
- 510 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and  
511 Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine*  
512 *Learning*, pp. 10764–10799. PMLR, 2023.  
513
- 514 Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Mari-  
515 beth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, et al. Improving alignment of  
516 dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- 517 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
518 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd  
519 of models. *arXiv preprint arXiv:2407.21783*, 2024.  
520
- 521 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,  
522 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the  
523 science of language models. *arXiv preprint arXiv:2402.00838*, 2024.
- 524 David Halliday, Robert Resnick, and Jearl Walker. *Fundamental-Physics*. Unknown.
- 525 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,  
526 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for  
527 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint*  
528 *arXiv:2402.14008*, 2024.  
529
- 530 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
531 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
532 *arXiv:2009.03300*, 2020.
- 533 Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers.  
534 *arXiv preprint arXiv:2212.10071*, 2022.
- 535 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
536 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.  
537
- 538 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-  
539 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*  
*arXiv:2410.21276*, 2024.

- 540 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bam-  
541 ford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
542 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 543
- 544 Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han, and Rui Zhang. When can llms actually correct  
545 their own mistakes? a critical survey of self-correction of llms. *Transactions of the Association*  
546 *for Computational Linguistics*, 12:1417–1440, 2024.
- 547 Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli,  
548 Shariq Iqbal, Colton Bishop, Rebecca Roelofs, et al. Training language models to self-correct via  
549 reinforcement learning. *arXiv preprint arXiv:2409.12917*, 2024.
- 550
- 551 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
552 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-  
553 tion for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:  
554 9459–9474, 2020.
- 555 Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. Evaluating mathe-  
556 matical reasoning of large language models: A focus on error identification and correction. *arXiv*  
557 *preprint arXiv:2406.00755*, 2024.
- 558
- 559 Zhenwen Liang, Dian Yu, Xiaoman Pan, Wenlin Yao, Qingkai Zeng, Xiangliang Zhang, and Dong  
560 Yu. Mint: Boosting generalization in mathematical reasoning via multi-view fine-tuning. *arXiv*  
561 *preprint arXiv:2307.07951*, 2023.
- 562 Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su.  
563 Deductive verification of chain-of-thought reasoning. *Advances in Neural Information Processing*  
564 *Systems*, 36, 2024.
- 565
- 566 Elita Lobo, Chirag Agarwal, and Himabindu Lakkaraju. On the impact of fine-tuning on chain-of-  
567 thought reasoning. *arXiv preprint arXiv:2411.15382*, 2024.
- 568
- 569 Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qing-  
570 wei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning  
571 for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023.
- 572
- 573 Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhui Chen. General-reasoner:  
574 Advancing llm reasoning across all domains, 2025. URL <https://arxiv.org/abs/2505.14652>.
- 575
- 576 Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The  
577 sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165.  
Elsevier, 1989.
- 578
- 579 Ning Miao, Yee Whye Teh, and Tom Rainforth. Selfcheck: Using llms to zero-shot check their own  
580 step-by-step reasoning. *arXiv preprint arXiv:2308.00436*, 2023.
- 581
- 582 Microsoft Research. Phi-3 technical report: A highly capable language model locally on your phone.  
583 *arXiv preprint arXiv:2404.14219*, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 584
- 585 Vishwas Mruthyunjaya, Pouya Pezeshkpour, Estevam Hruschka, and Nikita Bhutani. Rethinking  
586 language models as symbolic knowledge graphs. *arXiv preprint arXiv:2308.13676*, 2023.
- 587
- 588 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
589 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
low instructions with human feedback. *Advances in neural information processing systems*, 35:  
27730–27744, 2022.
- 590
- 591 D. C. Pandey. *IIT JEE Physics 35 years chapter wise solved papers*. Unknown.
- 592
- 593 Xinyu Pang, Ruixin Hong, Zhanke Zhou, Fangrui Lv, Xinwei Yang, Zhilong Liang, Bo Han, and  
Changshui Zhang. Physics reasoner: Knowledge-augmented reasoning for solving physics prob-  
lems with large language models, 2024. URL <https://arxiv.org/abs/2412.13791>.

- 594 Avinash Patil and Aryan Jadon. Advancing reasoning in large language models: Promising methods  
595 and approaches, 2025. URL <https://arxiv.org/abs/2502.03671>.
- 596
- 597 A. A. Pinsky. *Problems in Physics*. Mir Publishers, 1989.
- 598
- 599 PW Live. H.c. verma solutions and study materials for school prep and exams, 2024. URL <https://www.pw.live/school-prep/exams/hc-verma-solutions>.
- 600
- 601 Jianhao Shen, Yichun Yin, Lin Li, Lifeng Shang, Xin Jiang, Ming Zhang, and Qun Liu. Generate &  
602 rank: A multi-task framework for math word problems. *arXiv preprint arXiv:2109.03034*, 2021.
- 603
- 604 Hakim Sidahmed, Samrat Phatale, Alex Hutcheson, Zhuonan Lin, Zhang Chen, Zac Yu, Jarvis Jin,  
605 Simral Chaudhary, Roman Komarytsia, Christiane Ahlheim, et al. Parameter efficient reinforce-  
606 ment learning from human feedback. *arXiv preprint arXiv:2403.10704*, 2024.
- 607
- 608 Gaurav Srivastava, Shuxiang Cao, and Xuan Wang. Towards reasoning ability of small language  
609 models, 2025. URL <https://arxiv.org/abs/2502.11569>.
- 610
- 611 Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai  
612 Yu. Scieval: A multi-level large language model evaluation benchmark for scientific research.  
613 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19053–19061,  
2024.
- 614
- 615 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,  
616 Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal under-  
617 standing across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024a.
- 618
- 619 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-  
620 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma  
2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024b.
- 621
- 622 Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- 623
- 624 Paul A. Tipler. *Physics (Vol. I and II)*. W. H. Freeman and Company, 1999.
- 625
- 626 Gladys Tyen, Hassan Mansoor, Victor Carbune, Peter Chen, and Tony Mak. LLMs cannot find rea-  
627 soning errors, but can correct them given the error location. In Lun-Wei Ku, Andre Martins, and  
628 Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp.  
13894–13908, Bangkok, Thailand and virtual meeting, August 2024. Association for Computa-  
629 tional Linguistics. URL <https://aclanthology.org/2024.findings-acl.826>.
- 630
- 631 Tu Vu, Tong Wang, Tsendsuren Munkhdalai, Alessandro Sordani, Adam Trischler, Andrew  
632 Mattarella-Micke, Subhransu Maji, and Mohit Iyyer. Exploring and predicting transferability  
across nlp tasks. *arXiv preprint arXiv:2005.00770*, 2020.
- 633
- 634 Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi  
635 Song, Mingjie Zhan, and Hongsheng Li. Mathcoder: Seamless code integration in llms for en-  
hanced mathematical reasoning. *arXiv preprint arXiv:2310.03731*, 2023.
- 636
- 637 Shangshang Wang, Julian Asilis, Ömer Faruk Akgül, Enes Burak Bilgin, Ollie Liu, and Willie  
638 Neiswanger. Tina: Tiny reasoning models via lora. *arXiv preprint arXiv:2504.15777*, 2025.
- 639
- 640 Yifei Wang, Yuheng Chen, Wanting Wen, Yu Sheng, Linjing Li, and Daniel Dajun Zeng. Unveil-  
641 ing factual recall behaviors of large language models through knowledge neurons, 2024. URL  
<https://arxiv.org/abs/2408.03247>.
- 642
- 643 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-  
644 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language  
645 models. *arXiv preprint arXiv:2206.07682*, 2022a.
- 646
- 647 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in  
neural information processing systems*, 35:24824–24837, 2022b.

- 648 Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. Evaluating mathematical  
649 reasoning beyond accuracy. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
650 volume 39, pp. 27723–27730, 2025.
- 651
- 652 Yisheng Xiao, Lijun Wu, Junliang Guo, Juntao Li, Min Zhang, Tao Qin, and Tie-yan Liu. A survey  
653 on non-autoregressive generation for neural machine translation and beyond. *IEEE Transactions*  
654 *on Pattern Analysis and Machine Intelligence*, 45(10):11407–11427, 2023.
- 655
- 656 Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosenberg, Zhen Qin, Daniele Calandriello,  
657 Misha Khalman, Rishabh Joshi, Bilal Piot, Mohammad Saleh, Chi Jin, Tong Zhang, and Tianqi  
658 Liu. Building math agents with multi-turn iterative preference learning, 2025. URL <https://arxiv.org/abs/2409.02392>.
- 659
- 660 Huimin Xu, Xin Mao, Feng-Lin Li, Xiaobao Wu, Wang Chen, Wei Zhang, and Anh Tuan Luu.  
661 Full-step-dpo: Self-supervised preference optimization with step-wise rewards for mathematical  
662 reasoning. *arXiv preprint arXiv:2502.14356*, 2025.
- 663
- 664 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik  
665 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-*  
666 *vances in Neural Information Processing Systems*, 36, 2024.
- 667
- 668 Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong Pasupat, Jure Leskovec, Percy Liang,  
669 Ed H Chi, and Denny Zhou. Large language models as analogical reasoners. *arXiv preprint*  
670 *arXiv:2310.01714*, 2023.
- 671
- 672 Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou,  
673 and Jingren Zhou. Scaling relationship on learning mathematical reasoning with large language  
674 models, 2024. URL <https://openreview.net/forum?id=cij00f8u35>.
- 675
- 676 Yunxiang Zhang, Muhammad Khalifa, Lajanugen Logeswaran, Jaekyeom Kim, Moontae Lee,  
677 Honglak Lee, and Lu Wang. Small language models need strong verifiers to self-correct reason-  
678 ing. *arXiv preprint arXiv:2404.17140*, 2024.
- 679
- 680 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in  
681 large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- 682
- 683 Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia,  
684 Linqi Song, Mingjie Zhan, et al. Solving challenging math word problems using gpt-4 code  
685 interpreter with code-based self-verification. *arXiv preprint arXiv:2308.07921*, 2023.

## 685 A APPENDIX

### 686 A.1 PHYSICSQA

#### 687 CONSTRUCTION

688

689

690

691 The PhysicsQA dataset was created to address the gap between foundational knowledge bench-  
692 marks and challenging problems. It consists of 370 diverse, intermediate-level  
693 high school physics problems. These problems were sourced from standard Indian JEE preparation  
694 materials from 2000-2010. Each problem is paired with a correct Chain-of-Thought (CoT) solu-  
695 tion, which was verified using GPT-4 with human annotators in the loop. This detailed structure  
696 allows for a step-by-step evaluation of an LLM’s conceptual understanding and calculation abilities,  
697 moving beyond simple final answer accuracy. The problems are designed to be challenging, often  
698 requiring multiple concepts, intricate calculations, and multi-hop reasoning. A separate fine-tuning  
699 dataset, consisting of 2,494 high-school level physics questions also sourced from JEE materials,  
700 was curated for the warm-up and RL training stages of our methodology.

701 The meticulous five-step process used to curate these challenging problems is illustrated in the figure  
below.

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

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 4: Topic-wise Distribution in PhysicsQA

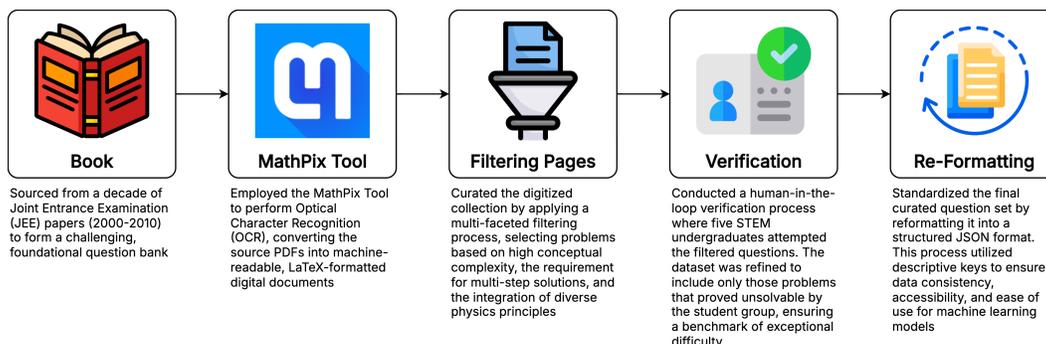


Figure 3: The dataset creation pipeline for PhysicsQA. The process begins with sourcing questions from Joint Entrance Examination (JEE) papers from 2000-2010. These are digitized into LaTeX format using the MathPix OCR tool. A rigorous filtering stage follows, where problems are selected for high conceptual complexity and the need for multi-step solutions. A human-in-the-loop verification process then ensures a high difficulty standard by selecting only those questions that a group of five STEM undergraduates could not solve. Finally, the curated problems are standardized into a structured JSON format with descriptive keys, making them suitable for machine learning applications.

DATASET SAMPLES

PhysicsQA

The following figures showcase the diversity of problems within the PhysicsQA dataset, spanning different physics domains and requiring a range of reasoning skills.

756 "id": "811",  
 757  
 758 "question": "Prove that the pressure of an ideal gas is numerically equal to two third of the mean  
 759 translational kinetic energy per unit volume of the gas.",  
 760  
 761 "cot\_solution": "  
 762 # Step 1: Understand the problem and identify what needs to be proved. We need to prove that the pressure  
 763 of an ideal gas is numerically equal to two-thirds of the mean translational kinetic energy per unit volume of  
 764 the gas.  
 765 # Step 2: Recall the formula for translational kinetic energy per unit volume. Translational kinetic energy per  
 766 unit volume (E) is given by the formula:  $E = \frac{1}{2}(\text{mass per unit volume})(\bar{v}^2)$  where  $\bar{v}^2$  is the mean squared  
 767 velocity of the gas molecules.  
 768  
 769 # Step 6: Rearrange the equation to solve for pressure. Rearranging the equation to solve for p, we get:  
 770  $p = \frac{2}{3}E$ . The final answer is:  $\left\{ \frac{2}{3}E \right\}$ ,  
 771  
 772 "subject": "Physics",  
 773  
 774 "source": "JEE Questions",  
 775  
 776 "chapter\_name": "Thermodynamics and heat transfer",  
 777  
 778 "class": "High School"

781 Figure 4: A sample from the PhysicsQA dataset involving a theoretical proof in thermodynamics.  
 782 This problem tests the model's ability to reason with abstract concepts and manipulate fundamental  
 783 formulas, such as the relationship between pressure, density, and kinetic energy, to derive a result  
 784 rather than calculating a numerical value. It highlights the dataset's inclusion of conceptual reason-  
 785 ing tasks.  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807  
 808  
 809

810 "id": "846",  
 811  
 812 "question": "Example 17: Uranium ores on the earth at the present time typically have a composition  
 813 consisting of 99.3% of the isotope  $\{{}_{92}U^{238}$  and 0.7% of the isotope  $\{{}_{92}U^{235}$ . The half-lives of these  
 814 isotopes are  $4.47 \times 10^9$  year and  $7.04 \times 10^8$  year, respectively. If these isotopes were equally abundant when  
 815 the earth was formed, estimate the age of the earth.",  
 816  
 817 "cot\_solution": "  
 818 # Step 1: Define the initial conditions and the given information.  
 819 We are given that uranium ores on the earth at the present time typically have a composition consisting of  
 820 99.3% of the isotope  $\{{}_{92}U^{238}$  and 0.7% of the isotope  $\{{}_{92}U^{235}$ . The half-lives of these isotopes are  
 821  $4.47 \times 10^9$  year and  $7.04 \times 10^8$  year respectively.  
 822  
 823 # Step 2: Determine the decay constants for both isotopes.  
 824 The decay constant  $\lambda$  is related to the half-life  $t_{\frac{1}{2}}$  by the equation  $\lambda = \frac{0.693}{t_{\frac{1}{2}}}$ . Therefore, the decay constants  
 825 for the two isotopes are  $\lambda_1 = \frac{0.693}{4.47 \times 10^9}$  and  $\lambda_2 = \frac{0.693}{7.04 \times 10^8}$ .  
 826  
 827 ...  
 828  
 829 # Step 6: Substitute the values of the decay constants and calculate the age of the earth.  
 830 Substituting the values of the decay constants, we have  $t = \frac{1}{\frac{0.693}{7.04 \times 10^8} - \frac{0.693}{4.47 \times 10^9}} \ln\left(\frac{99.3}{0.7}\right)$ . Evaluating this  
 831 expression, we get  $t = 5.97 \times 10^9$  year.  
 832  
 833 The final answer is:  $\{5.97 \times 10^9\}$ ,  
 834  
 835 "subject": "Physics",  
 836  
 837 "source": "Optics and modern physics",  
 838  
 839 "chapter\_name": "Nuclear physics and Radioactivity",  
 840  
 841 "class": "12th"

842  
 843 Figure 5: A sample from the nuclear physics and radioactivity chapter. This problem requires  
 844 multi-step numerical calculation involving half-life, decay constants, and logarithms. It assesses the  
 845 model's precision with scientific notation and its ability to apply formulas to a real-world scenario  
 846 (estimating the age of the Earth), demonstrating the dataset's coverage of advanced topics and com-  
 847 plex calculations.  
 848  
 849  
 850  
 851  
 852  
 853  
 854  
 855  
 856  
 857  
 858

## 859 JEEBench

860  
 861 JEEBench, derived from the Joint Entrance Examination (JEE) Advanced papers in India, repre-  
 862 sents a benchmark of exceptionally high difficulty, designed to challenge the top echelon of human  
 863 students. The problems are renowned for requiring deep conceptual understanding, mathematical  
 fluency, and the ability to synthesize multiple topics under pressure.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

"description": "JEE Adv 2016 Paper 1",

"index": 2,

"subject": "phy",

"type": "MCQ",

"question": "A uniform wooden stick of mass  $1.6\{\sim kg\}$  and length  $l$  rests in an inclined manner on a smooth, vertical wall of height  $h(< l)$  such that a small portion of the stick extends beyond the wall. The reaction force of the wall on the stick is perpendicular to the stick. The stick makes an angle of  $30^\circ$  with the wall and the bottom of the stick is on a rough floor. The reaction of the wall on the stick is equal in magnitude to the reaction of the floor on the stick. The ratio  $\frac{h}{l}$  and the frictional force  $f$  at the bottom of the stick are

$(g = 10\{\sim ms\}\{\sim s\}^2)$

(A)  $\frac{h}{l} = \frac{\sqrt{3}}{16}, f = \frac{16\sqrt{3}}{3}\{\sim N\}$

(B)  $\frac{h}{l} = \frac{3}{16}, f = \frac{16\sqrt{3}}{3}\{\sim N\}$

(C)  $\frac{h}{l} = \frac{3\sqrt{3}}{16}, f = \frac{8\sqrt{3}}{3}\{\sim N\}$

(D)  $\frac{h}{l} = \frac{3\sqrt{3}}{16}, f = \frac{16\sqrt{3}}{3}\{\sim N\}$ ,

"gold": "D"

Figure 6: A high-complexity problem in static equilibrium. This question requires a model to construct a complete physical model (a free-body diagram), apply both force ( $\sum \vec{F} = 0$ ) and torque ( $\sum \vec{\tau} = 0$ ) equilibrium conditions, and solve a resulting system of simultaneous equations involving trigonometry. It is a rigorous test of translating a physical setup into a solvable mathematical framework.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

"description": "JEE Adv 2016 Paper 1",  
 "index": 12,  
 "subject": "phy",  
 "type": "MCQ(multiple)",  
 "question": "The position vector  $\vec{r}$  of a particle of mass  $m$  is given by the following equation  

$$\vec{r}(t) = \alpha t^3 \hat{i} + \beta t^2 \hat{j}$$
 where  $\alpha = \frac{10}{3} \{ \sim m \} \{ \sim s \}^{-3}$ ,  $\beta = 5 \{ \sim m \} \{ \sim s \}^{-2}$  and  $m = 0.1 \{ \sim kg \}$ . At  $t = 1 \{ \sim s \}$ , which of the following statement(s) is(are) true about the particle?  
 (A) The velocity  $\vec{v}$  is given by  $\vec{v} = (10\hat{i} + 10\hat{j}) \{ ms \}^{-1}$   
 (B) The angular momentum  $\vec{L}$  with respect to the origin is given by  $\vec{L} = -\left(\frac{5}{3}\right) \hat{k} \{ \sim N \} \{ \sim m \}$   
 (C) The force  $\vec{F}$  is given by  $\vec{F} = (\hat{i} + 2\hat{j}) \{ N \}$   
 (D) The torque  $\vec{\tau}$  with respect to the origin is given by  $\vec{\tau} = -\left(\frac{20}{3}\right) \hat{k} \{ \sim N \} \{ \sim m \}$ ",  
 "gold": "ABD"

Figure 7: A calculus-intensive vector mechanics problem presented in a multiple-correct format. A model must perform several distinct vector operations—differentiation to find velocity and acceleration, and cross products to find angular momentum ( $\vec{L} = \vec{r} \times m\vec{v}$ ) and torque ( $\vec{\tau} = \vec{r} \times \vec{F}$ ). The multiple-correct nature demands that each option be independently and accurately verified, testing procedural stamina.

"description": "JEE Adv 2017 Paper 2",  
 "index": 1,  
 "subject": "phy",  
 "type": "MCQ",  
 "question": "Consider an expanding sphere of instantaneous radius  $R$  whose total mass remains constant. The expansion is such that the instantaneous density  $\rho$  remains uniform throughout the volume. The rate of fractional change in density  $\left(\frac{1}{\rho} \frac{d\rho}{dt}\right)$  is constant. The velocity  $v$  of any point on the surface of the expanding sphere is proportional to  
 [A]  $R$   
 [B]  $R^3$   
 [C]  $\frac{1}{R}$   
 [D]  $R^{\frac{2}{3}}$ ",  
 "gold": "A"

Figure 8: An abstract problem requiring **differential reasoning**. Unlike applying a known formula, this question demands the ability to derive a relationship from first principles. The model must relate density to radius, differentiate their relationship with respect to time, and interpret the result to find the proportionality, testing a model's ability to reason with changing quantities.

## COMPARATIVE ANALYSIS: JEEBENCH VS. PHYSICSQA

While JeeBench exemplifies a high level of difficulty in computational and integrative physics, the PhysicsQA dataset is deliberately designed to be a more robust and insightful benchmark for evaluating the reasoning of large language models for several key reasons:

1. **Process vs. Outcome:** JEEBench, with its multiple-choice format, exclusively evaluates the final answer. A correct answer could be achieved through flawed reasoning or even a lucky guess. In contrast, PhysicsQA’s structured **chain-of-thought (CoT) solutions** demand the generation of a complete, step-by-step reasoning path. This makes it a superior diagnostic tool, as it assesses the validity of the problem-solving process itself, not just the outcome.
2. **Diversity of Reasoning Tasks:** The JEEBench samples, while complex, are primarily focused on intricate calculation and synthesis within classical mechanics. PhysicsQA demonstrates broader cognitive diversity by including tasks such as **theoretical proofs** (as seen in Figure 4) and **real-world estimation problems** (Figure 5). This provides a more holistic evaluation of scientific intelligence beyond competition-style problems.
3. **Generation vs. Discrimination:** The fundamental task in JeeBench is discriminative selecting the correct option(s) from a provided list. PhysicsQA requires a generative approach creating a detailed solution from scratch. For an AI model, generating a coherent, logically sound, and mathematically correct multi-step solution is a significantly more demanding and authentic measure of true problem-solving ability than choosing an answer.

While the difficulty of JEEBench is undeniable, the architectural design of the PhysicsQA dataset with its focus on the reasoning process, diverse task types, and generative nature—makes it a more powerful and precise instrument for measuring and advancing the scientific reasoning capabilities of large language models.

**SciEval-Static**

The SciEval-Static dataset provides a benchmark for evaluating models on foundational physics problems. The samples are typically multiple-choice questions that require direct application of core principles and formulas, assessing both conceptual knowledge and computational accuracy.

```

"question": "If a 2kg object is constantly accelerated from  $0 \frac{m}{s}$  to  $16 \frac{m}{s}$  over 6 s, how much power must be
applied at  $t = 1$ ?
  A.  $16/3$  N
  B.  $8/3$  N
  C.  $32/3$  N
  D.  $4$  N",
"answer": ["A"],
"category": "physics",
"topic": "Work and Energy",
"ability": "Scientific Calculation",
"type": "multiple-choice",
"task_name": "SocraticQA",
"id": "415"

```

Figure 9: A procedural problem from ‘Work and Energy’ that tests **sequential calculation**. The model must first determine the object’s acceleration from kinematic data and then use that result to calculate the force and instantaneous power. This question assesses the ability to follow a clear, step-by-step computational procedure using fundamental formulas ( $a = \Delta v / \Delta t$ ,  $F = ma$ ,  $P = Fv$ ).

1026  
 1027 "question": "If a concave and convex lens of equal focus are kept in contact, what is the effective focal length?  
 1028 A. Infinite focal length (equivalent to plane glass)  
 1029 B. Zero focal length  
 1030 C. Equal to the original focal length  
 1031 D. Half of the original focal length",  
 1032 "answer": ["A"],  
 1033 "category": "physics",  
 1034 "topic": "Refraction",  
 1035 "ability": "Scientific Calculation",  
 1036 "type": "multiple-choice",  
 1037 "task\_name": "SocraticQA",  
 1038 "id": "809"

1043  
 1044  
 1045 Figure 10: A **conceptual reasoning** question from optics testing knowledge of the lens combination  
 1046 formula and, critically, scientific sign conventions. The solution hinges on recognizing that a convex  
 1047 ( $+f$ ) and a concave ( $-f$ ) lens of equal focal length mathematically negate each other. This problem  
 1048 highlights the dataset’s emphasis on testing core physics principles rather than complex arithmetic.

1049  
 1050  
 1051 "question": "If a projectile is shot at a velocity of  $7\frac{m}{s}$  and an angle of  $\frac{\pi}{3}$ , how far will the projectile travel  
 1052 before landing?  
 1053 A. 5.12 m  
 1054 B. 4.33 m  
 1055 C. 6.28 m  
 1056 D. 3.67 m",  
 1057 "answer": ["B"],  
 1058 "category": "physics",  
 1059 "topic": "2D Motion",  
 1060 "ability": "Scientific Calculation",  
 1061 "type": "multiple-choice",  
 1062 "task\_name": "SocraticQA",  
 1063 "id": "97"

1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071 Figure 11: A standard projectile motion problem that evaluates a model’s ability to perform a **single-**  
 1072 **step scientific calculation**. Success requires identifying the correct formula for horizontal range  
 1073 ( $R = \frac{v_0^2 \sin(2\theta)}{g}$ ) and accurately substituting the given values. This sample is representative of the  
 1074 dataset’s inclusion of foundational, formula-driven problems.

1075  
 1076 **MMLU-HighSchool**  
 1077

1078 The MMLU-HighSchool dataset contains questions that primarily evaluate a model’s conceptual  
 1079 grasp of physics. The problems often require qualitative reasoning and an understanding of founda-  
 tional principles rather than intricate multi-step calculations.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090

```
"question": "The plates of a capacitor are charged to a potential difference of 5 V. If the capacitance is 2 mF,
what is the charge on the positive plate?",

"input": "
  A. 0.005 °C,
  B. 0.01 °C,
  C. 0.02 °C,
  D. 0.5 °C",

"answer": "B",

"topic": "Electromagnetism"
```

1091  
1092  
1093  
1094  
1095  
1096

Figure 12: A foundational problem in electromagnetism testing direct **knowledge recall** and application of the core capacitor equation ( $Q = CV$ ). This question also assesses a model’s ability to correctly interpret scientific prefixes (mF for milli-Farads), a crucial step in ensuring accurate numerical computation in a straightforward context.

1097  
1098

1099  
1100  
1101  
1102

```
"question": "It is known that a lab cart is moving east at 25 cm/s at time t1 = 0.10 s, and then moving east at
15 cm/s at t2 = 0.20 s. Is this enough information to determine the direction of the net force acting on the
cart between t1 and t2?",

"input": "
  A. Yes, since we know the cart is slowing down, its momentum change is opposite the direction of
  movement, and the net force is in the direction of momentum change.,
  B. No, because we don't know whether forces such as friction or air resistance might be acting on the cart.,
  C. No, because we don't know the mass of the cart.,
  D. Yes, since we know the cart keeps moving to the east, the net force must be in the direction of motion.",

"answer": "A",

"topic": "Mechanics"
```

1111  
1112  
1113  
1114  
1115  
1116  
1117

Figure 13: A question focused on **qualitative reasoning** in mechanics. It requires no calculation but instead tests a deep conceptual understanding of Newton’s Second Law ( $\vec{F}_{net} = m\vec{a}$ ). The model must deduce the direction of the net force from the change in velocity (acceleration), correctly identifying that a decelerating object has a net force opposing its motion.

1118  
1119

1120  
1121  
1122

```
"question": "A sound wave with frequency f travels through air at speed v. With what speed will a sound wave
with frequency 4f travel through the air?",

"input": "A. v/4, B. v, C. 2v, D. 4v",

"answer": "B",

"topic": "Waves"
```

1123  
1124  
1125  
1126  
1127  
1128  
1129

1130  
1131  
1132  
1133

Figure 14: A **conceptual pitfall** question from the topic of waves. It is designed to test a model’s understanding of a core principle: the speed of a wave is determined by the properties of its medium, not its frequency. Success requires resisting the common misconception derived from the wave equation ( $v = f\lambda$ ), highlighting the dataset’s role in evaluating robustness against common distractors.

1134 **MMLU-College**

1135  
1136 The MMLU-College dataset challenges models with problems that reflect the complexity and in-  
1137 terdisciplinary nature of a university-level physics education. These questions often require the  
1138 synthesis of multiple concepts and a high degree of factual precision.

1139  
1140  
1141  
1142  
1143

1144 "question": "The quantum efficiency of a photon detector is 0.1. If 100 photons are sent into the detector, one  
1145 after the other, the detector will detect photon.",

1146 "input": "

- 1147     A. an average of 10 times, with an rms deviation of about 4,  
1148     B. an average of 10 times, with an rms deviation of about 3,  
1149     C. an average of 10 times, with an rms deviation of about 1,  
1150     D. an average of 10 times, with an rms deviation of about 0.1",

1151 "answer": "B",

1152  
1153 "topic": "Modern\_Physics"

1154  
1155

1156 Figure 15: An **interdisciplinary reasoning** problem that merges a concept from modern physics  
1157 (quantum efficiency) with probability theory. The solution requires not just calculating the expected  
1158 value but also understanding statistical fluctuations by applying the formula for the standard devi-  
1159 ation of a binomial distribution ( $\sigma = \sqrt{np(1-p)}$ ). This sample showcases problems that test the  
1160 mathematical rigor expected in higher education.

1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170

1171 "question": "The coefficient of static friction between a small coin and the surface of a turntable is 0.30. The  
1172 turntable rotates at 33.3 revolutions per minute. What is the maximum distance from the center of the  
1173 turntable at which the coin will not slide?",

1174 "input": "

- 1175     A. 0.024 m,  
1176     B. 0.048 m,  
1177     C. 0.121 m,  
1178     D. 0.242 m",

1179 "answer": "D",

1180  
1181 "topic": "Mechanics"

1182  
1183  
1184  
1185  
1186  
1187

1184 Figure 16: A classic university-level mechanics problem that requires the **synthesis of multiple**  
1185 **concepts**: circular motion (centripetal force) and static friction. A critical component is the **proce-**  
1186 **dural accuracy** needed for unit conversion—transforming revolutions per minute into radians per  
1187 second. This problem tests a model’s ability to construct a solution by integrating different physical  
principles.

```

1188
1189 "question": "Electromagnetic radiation provides a means to probe aspects of the physical universe. Which of
1190 the following statements regarding radiation spectra is NOT correct?",
1191
1192 "input": "
1193     A. Lines in the infrared, visible, and ultraviolet regions of the spectrum reveal primarily the nuclear
1194     structure of the sample.,
1195     B. The wavelengths identified in an absorption spectrum of an element are among those in its emission
1196     spectrum.,
1197     C. Absorption spectra can be used to determine which elements are present in distant stars.,
1198     D. Spectral analysis can be used to identify the composition of galactic dust.",
1199
1200 "answer": "A",
1201
1202 "topic": "Electromagnetism"

```

Figure 17: A knowledge-intensive question requiring **factual precision** regarding electromagnetic spectra. The 'NOT correct' format tests a model's ability to critically evaluate several statements and identify the specific scientific inaccuracy—in this case, confusing electronic energy level transitions with nuclear structure. This question assesses the breadth and depth of a model's domain-specific knowledge base.

## A.2 RETRIEVAL-AUGMENTED GENERATION (RAG) FOR FORMULAIC ACCURACY

To enhance the model's performance on quantitative problems and mitigate the risk of hallucinating incorrect formulas, we implemented a Retrieval-Augmented Generation (RAG) system. It is crucial to clarify the nature of the retrieval corpus to accurately represent the task posed to the LLM.

Contrary to a system that retrieves from a full textbook, our RAG corpus is intentionally constrained to a concise formula sheet. This knowledge base contains only essential information such as physical constants (e.g., speed of light, Planck's constant), fundamental laws (e.g., Newton's Laws), and topic-specific equations (e.g., for projectile motion or LCR circuits). The full document provided to the RAG system is detailed in the appendix.

This design choice is deliberate. By providing access only to the "tools" of physics the formulas themselves we ensure the model cannot simply find and replicate a pre-existing worked-out solution. The LLM is still required to perform all the critical reasoning steps: understanding the problem, selecting the appropriate formulas from the retrieved context, manipulating them correctly, and executing the final calculation. This approach tests the model's ability to *apply* knowledge, not merely to find and copy it.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

**5 Electricity and Magnetism**

**5.1: Electrostatics**

Coulomb's law:  $\vec{F} = \frac{1}{4\pi\epsilon_0} \frac{q_1 q_2}{r^2} \hat{r}$

Electric field:  $\vec{E}(\vec{r}) = \frac{1}{4\pi\epsilon_0} \frac{q}{r^2} \hat{r}$

Electrostatic energy:  $U = \frac{1}{4\pi\epsilon_0} \frac{q_1 q_2}{r}$

Electrostatic potential:  $V = \frac{1}{4\pi\epsilon_0} \frac{q}{r}$

$dV = -\vec{E} \cdot d\vec{r}$ ,  $V(\vec{r}) = -\int_{\infty}^{\vec{r}} \vec{E} \cdot d\vec{r}$

Electric dipole moment:  $\vec{p} = q\vec{d}$

Potential of a dipole:  $V = \frac{1}{4\pi\epsilon_0} \frac{p \cos \theta}{r^2}$

Force between plates of a parallel plate capacitor:  $F = \frac{Q^2}{2\epsilon_0 A}$

Energy stored in capacitor:  $U = \frac{1}{2} CV^2 = \frac{Q^2}{2C} = \frac{1}{2} QV$

Energy density in electric field  $E$ :  $U/V = \frac{1}{2} \epsilon_0 E^2$

Capacitor with dielectric:  $C = \epsilon_0 \epsilon_r \frac{A}{d}$

**5.2: Gauss's Law and its Applications**

Electric flux:  $\phi = \oint \vec{E} \cdot d\vec{S}$

Gauss's law:  $\oint \vec{E} \cdot d\vec{S} = q_{encl}/\epsilon_0$

Field of a uniformly charged ring on its axis

$E_x = \frac{1}{4\pi\epsilon_0} \frac{2\pi R q}{(\sigma^2 + z^2)^{3/2}}$

$E$  and  $V$  of a uniformly charged sphere:

$E = \begin{cases} \frac{1}{4\pi\epsilon_0} \frac{Qr}{R^3}, & \text{for } r < R \\ \frac{1}{4\pi\epsilon_0} \frac{Q}{r^2}, & \text{for } r \geq R \end{cases}$

$V = \begin{cases} \frac{Q}{4\pi\epsilon_0 R} \left( \frac{3}{2} - \frac{r^2}{2R^2} \right), & \text{for } r < R \\ \frac{Q}{4\pi\epsilon_0 r}, & \text{for } r \geq R \end{cases}$

$E$  and  $V$  of a uniformly charged spherical shell:

$E = \begin{cases} 0, & \text{for } r < R \\ \frac{1}{4\pi\epsilon_0} \frac{Q}{r^2}, & \text{for } r \geq R \end{cases}$

$V = \begin{cases} \frac{Q}{4\pi\epsilon_0 R}, & \text{for } r < R \\ \frac{Q}{4\pi\epsilon_0 r}, & \text{for } r \geq R \end{cases}$

Field of a line charge:  $E = \frac{\lambda}{2\pi\epsilon_0 r}$

Field of an infinite sheet:  $E = \frac{\sigma}{2\epsilon_0}$

Field in the vicinity of conducting surface:  $E = \frac{\sigma}{\epsilon_0}$

**5.3: Capacitors**

Capacitance:  $C = q/V$

Parallel plate capacitor:  $C = \epsilon_0 A/d$

Spherical capacitor:  $C = \frac{4\pi\epsilon_0 a b}{b-a}$

Cylindrical capacitor:  $C = \frac{2\pi\epsilon_0 l}{\ln(r_2/r_1)}$

Capacitors in parallel:  $C_{eq} = C_1 + C_2$

Capacitors in series:  $\frac{1}{C_{eq}} = \frac{1}{C_1} + \frac{1}{C_2}$

Force between plates of a parallel plate capacitor:  $F = \frac{Q^2}{2\epsilon_0 A}$

Energy stored in capacitor:  $U = \frac{1}{2} CV^2 = \frac{Q^2}{2C} = \frac{1}{2} QV$

Energy density in electric field  $E$ :  $U/V = \frac{1}{2} \epsilon_0 E^2$

Capacitor with dielectric:  $C = \epsilon_0 \epsilon_r \frac{A}{d}$

**5.4: Current electricity**

Current density:  $\vec{j} = i/A = \sigma \vec{E}$

Drift speed:  $v_d = \frac{1}{nA} \frac{dQ}{dt} = \frac{I}{nA}$

Resistance of a wire:  $R = \rho l/A$ , where  $\rho = 1/\sigma$

Temp. dependence of resistance:  $R = R_0(1 + \alpha \Delta T)$

Ohm's law:  $V = iR$

Kirchhoff's Laws: (i) The Junction Law: The algebraic sum of all the currents directed towards a node is zero i.e.,  $\sum_{node} I_i = 0$ . (ii) The Loop Law: The algebraic sum of all the potential differences along a closed loop in a circuit is zero i.e.,  $\sum_{loop} \Delta V_i = 0$ .

Resistors in parallel:  $\frac{1}{R_{eq}} = \frac{1}{R_1} + \frac{1}{R_2}$

Resistors in series:  $R_{eq} = R_1 + R_2$

Wheatstone bridge:

Balanced if  $R_1/R_2 = R_3/R_4$ .

Electric Power:  $P = V^2/R = IR = IV$

**5.5: Magnetism**

Galvanometer as an Ammeter:  $i_g G = (i - i_g) S$

Galvanometer as a Voltmeter:  $V_{AB} = i_g (R + G)$

Charging of capacitors:  $q(t) = CV [1 - e^{-t/\tau}]$

Discharging of capacitors:  $q(t) = q_0 e^{-t/\tau}$

Time constant in RC circuit:  $\tau = RC$

Peltier effect:  $\text{emf } e = \frac{\Delta T}{\Delta \theta} = \frac{\text{Peltier heat}}{\text{charge transferred}}$

Seebeck effect:  $e = \frac{\Delta T}{\Delta \theta} = \frac{\text{Thomson heat}}{\text{charge transferred}} = \sigma \Delta T$

Thomson effect:  $\text{emf } e = \frac{\Delta T}{\Delta \theta} = \frac{\text{Thomson heat}}{\text{charge transferred}} = \sigma \Delta T$

Faraday's law of electrolysis: The mass deposited is  $m = Zit = \frac{1}{2} Eit$

where  $i$  is current,  $t$  is time,  $Z$  is electrochemical equivalent,  $E$  is chemical equivalent, and  $F = 96485 \text{ C/g}$  is Faraday constant.

**5.6: Magnetic Field due to Current**

Energy of a magnetic dipole placed in  $\vec{B}$ :  $U = -\vec{\mu} \cdot \vec{B}$

Hall effect:  $V_H = \frac{B i}{n e d}$

Biot-Savart law:  $d\vec{B} = \frac{\mu_0}{4\pi} \frac{i d\vec{l} \times \vec{r}}{r^3}$

Field due to a straight conductor:  $B = \frac{\mu_0 i}{2\pi r} (\cos \theta_1 - \cos \theta_2)$

Field due to an infinite straight wire:  $B = \frac{\mu_0 i}{2\pi r}$

Force between parallel wires:  $\frac{dF}{dl} = \frac{\mu_0 i_1 i_2}{2\pi r}$

Field on the axis of a ring:  $B = \frac{\mu_0 i}{2} \frac{R^2}{(R^2 + z^2)^{3/2}}$

Field at the centre of an arc:  $B = \frac{\mu_0 i \theta}{4\pi R}$

Field at the centre of a ring:  $B = \frac{\mu_0 i}{2R}$

Ampere's law:  $\oint \vec{B} \cdot d\vec{l} = \mu_0 I_{encl}$

Field inside a solenoid:  $B = \mu_0 n i$ ,  $n = \frac{N}{l}$

Field inside a toroid:  $B = \frac{\mu_0 N i}{2\pi r}$

Field of a bar magnet:  $B_1 = \frac{\mu_0 M}{4\pi r^3}$ ,  $B_2 = \frac{\mu_0 M}{2\pi r^3}$

Angle of dip:  $B_1 = B \cos \delta$

Tangent galvanometer:  $B_H \tan \theta = \frac{\mu_0 N i}{2r}$ ,  $i = K \tan \theta$

Moving coil galvanometer:  $nAB = k\theta$ ,  $i = \frac{k}{nAB} \theta$

Time period of magnetometer:  $T = 2\pi \sqrt{\frac{I}{MBH}}$

Permeability:  $\vec{B} = \mu \vec{H}$

Figure 18: A representative snapshot of the formula sheet used as the knowledge base for our RAG system. The corpus is strictly limited to formulas, definitions, and physical constants, such as Coulomb's Law, capacitor equations, and the Biot-Savart law. The document explicitly **lacks** narrative explanations, problem-solving strategies, or worked-out examples. This design ensures that the RAG acts as a reliable formula reference, compelling the LLM to perform the core reasoning and problem-solving tasks independently.

## A.2.1 TRAINING DATASET SAMPLE

```

1296 "id": "721",
1297
1298 "question": "A car moves with a constant tangential acceleration  $a_t = \frac{2}{3} \frac{\text{m}}{\text{s}^2}$  along a horizontal surface
1299
1300 circumscribing a circle of radius  $r = 40\text{-m}$ . The coefficient of sliding friction between the wheels of the car
1301 and the surface is  $\mu_s = 0.20$ . What distance will the car ride without sliding if at the initial moment of time
1302 its velocity is equal to zero (take  $g = 10 \frac{\text{m}}{\text{s}^2}$  )",
1303
1304 "cot_solution": "
1305
1306 # Step 1: Identify the forces acting on the car and their components.
1307 The car is moving with both radial and tangential acceleration, provided by the frictional force. The frictional
1308 force  $f$  has radial ( $f_r$ ) and tangential ( $f_t$ ) components.
1309
1310 # Step 2: Express the radial and tangential components of the frictional force.
1311  $f_r = \frac{mv^2}{r}$  and  $f_t = ma_t$ .
1312
1313 # Step 3: State the condition to avoid sliding.
1314 To avoid sliding,  $f = \sqrt{f_r^2 + f_t^2} \leq \mu_s N$ , where  $N = mg$ .
1315
1316 # Step 4: Determine the point of sliding.
1317 At the sliding point,  $f_r^2 + f_t^2 = \mu_s^2 m^2 g^2$ .
1318
1319 # Step 5: Substitute the expressions for  $f_r$  and  $f_t$  into the equation for the sliding point.
1320 Substituting for  $f_r$  and  $f_t$ , we get:  $V^4 = (\mu_s^2 g^2 - a_t^2) r^2$ .
1321
1322 # Step 6: Use the equation of motion to relate the velocity to the distance traveled.
1323 For zero initial velocity, we have  $v^2 = 2a_t S$ .
1324
1325 # Step 7: Solve for the distance traveled by the car without sliding.
1326  $\Rightarrow S = \sqrt{\frac{(\mu_s^2 g^2 - a_t^2) r^2}{4a_t^2}} = 40\sqrt{2}\text{-m}$ .
1327
1328 The final answer is:  $\{40\sqrt{2}\}$ ",
1329
1330 "subject": "Physics",
1331
1332 "source": "Physics Module",
1333
1334 "chapter_name": "Mechanics",
1335
1336 "class": "11th"

```

Figure 19: A sample from our physics Training Dataset, mechanics that requires multi-hop reasoning. The solution involves integrating concepts from dynamics (frictional and centripetal forces) with kinematics to find the final distance. This type of problem tests the model’s ability to construct a coherent reasoning chain by combining multiple distinct physics principles, a key challenge in complex problem-solving.

## PRELIMINARY EXPLANATION

## 1. LARGE LANGUAGE MODELS

The rapid advancement of artificial intelligence, particularly with the development of Large Language Models (LLMs) built on the transformer architecture, has redefined the capabilities of natural language processing. These models now exhibit remarkable performance across various language-related tasks, such as text generation, question answering, translation, and summarization, often

1350 rivaling human-like comprehension. More intriguingly, LLMs have demonstrated emergent abilities  
1351 extending beyond their core functions, showing proficiency in tasks like commonsense reasoning,  
1352 code generation, and arithmetic. Several key factors have driven the evolution of LLMs, most no-  
1353 tably the exponential growth in available data and computational resources. Indeed, on the one  
1354 hand, social media platforms, digital libraries, and other sources have provided vast amounts of tex-  
1355 tual and multimedia information, enabling LLMs to be trained on extensive and diverse datasets. On  
1356 the other hand, the availability of powerful GPUs, TPUs, and distributed computing frameworks has  
1357 made it feasible to train models with billions, and even trillions, of parameters. Together, these two  
1358 factors have led LLMs to capture nuanced linguistic patterns, cultural context, and domain-specific  
1359 knowledge, enhancing their ability to generate coherent, contextually appropriate, and highly versa-  
1360 tile outputs.

## 1361 2. PROMPTING TECHNIQUES

1363 **Zero-shot :** Zero-shot prompting refers to a method where a large language model is given a  
1364 task purely through natural language instructions, without any example inputs or outputs. This  
1365 technique relies entirely on the model’s pretrained knowledge and generalization capabilities. The  
1366 primary need for zero-shot prompting arises in scenarios where labeled data is unavailable or when  
1367 rapid deployment is required without fine-tuning. Despite its simplicity, it has proven surprisingly  
1368 effective in many tasks like classification, summarization, and translation, particularly with larger  
1369 models like GPT-4. However, its performance can be inconsistent, especially for functions requiring  
1370 complex reasoning or domain-specific understanding. Furthermore, the outcome is highly sensitive  
1371 to prompt phrasing, and the lack of examples limits the model’s ability to grasp nuanced task  
1372 expectations.

1374 **Few-shot :** Few-shot prompting builds upon zero-shot by including a small number of input-output  
1375 examples within the prompt, helping the model better understand the task format and desired  
1376 response style. This technique addresses the limitations of zero-shot by offering in-context learning,  
1377 which can significantly enhance performance on structured tasks such as question answering or code  
1378 generation. It is especially valuable when collecting large training datasets, which are impractical,  
1379 but a few representative examples are available. Few-shot prompting has demonstrated notable  
1380 effectiveness, particularly when combined with chain-of-thought examples for reasoning-heavy  
1381 tasks. Nonetheless, it comes with limitations: the model’s performance may degrade if examples  
1382 are poorly chosen, it has limited capacity to retain many examples due to prompt length constraints,  
1383 and it still lacks persistent learning across sessions.

1384 **Chain-of-Thought (CoT):** Chain-of-Thought (CoT) prompting is an enhanced strategy developed  
1385 to augment the performance of large language models (LLMs) on complex reasoning tasks such as  
1386 arithmetic, commonsense, and symbolic reasoning. This method integrates intermediate reasoning  
1387 steps within the prompts, providing a more structured path towards the solution. The complexity of  
1388 reasoning steps is the most critical factor for the performance of LLMs on complex reasoning tasks.  
1389 Complexity-based prompting can be further enhanced by using the output selection method called  
1390 Complexity-based Consistency, alleviating the possibility that the model can take shortcuts during  
1391 reasoning.

## 1392 1393 3. FINE-TUNING [LORA TECHNIQUE]

1395 In the realm of language models, fine-tuning an existing language model to perform a specific task  
1396 on specific data is a common practice. This involves adding a task-specific head, if necessary, and  
1397 updating the weights of the neural network through backpropagation during the training process. It  
1398 is essential to note the distinction between this fine-tuning process and training from scratch. In the  
1399 latter scenario, the model’s weights are randomly initialized, while in fine-tuning, the weights are  
1400 already optimized to a certain extent during the pre-training phase. The decision of which weights  
1401 to optimize or update, and which ones to keep frozen, depends on the chosen technique. Fortunately,  
1402 parameter-efficient approaches for fine-tuning exist that have proven to be effective. Although most  
1403 such approaches have yielded less performance, Low Rank Adaptation (LoRA) has bucked this trend  
by even outperforming full fine-tuning in some cases, as a consequence of avoiding catastrophic

1404 forgetting (a phenomenon that occurs when the knowledge of the pretrained model is lost during the  
1405 fine-tuning process). LoRA is an improved fine-tuning method. Instead of fine-tuning all the weights  
1406 that constitute the weight matrix of the pre-trained large language model, two smaller matrices that  
1407 approximate this larger matrix are fine-tuned. These matrices constitute the LoRA adapter. This  
1408 fine-tuned adapter is then loaded into the pretrained model and used for inference.

#### 1409 1410 4. RAG

1411 Large language models (LLMs) have achieved remarkable success. However, they still face significant  
1412 limitations, especially in domain-specific or knowledge-intensive tasks, notably producing  
1413 "hallucinations" when handling queries beyond their training data or requiring current information.  
1414 To overcome challenges, Retrieval-Augmented Generation (RAG) enhances LLMs by retrieving relevant  
1415 document chunks from an external knowledge base through semantic similarity calculation.  
1416 By referencing external knowledge, RAG effectively reduces the problem of generating factually  
1417 incorrect content. Its integration into LLMs has resulted in widespread adoption, establishing RAG  
1418 as a key technology in advancing chat-bots and enhancing the suitability of LLMs for real-world  
1419 applications. RAG research shifted towards providing better information for LLMs to answer more  
1420 complex and knowledge-intensive tasks during the inference stage, leading to rapid development in  
1421 RAG studies. As research progressed, the enhancement of RAG was no longer limited to the inference  
1422 stage but began to incorporate more with LLM fine-tuning techniques. Among the optimization  
1423 methods for LLMs, RAG is often compared with Fine-tuning (FT) and prompt engineering. Prompt  
1424 engineering leverages a model's inherent capabilities with minimal necessity for external knowledge  
1425 and model adaptation. RAG can be likened to providing a model with a tailored textbook for  
1426 information retrieval, which is ideal for precise information retrieval tasks. In contrast, FT is comparable  
1427 to a student internalizing knowledge over time, suitable for scenarios requiring replication  
1428 of specific structures, styles, or formats. RAG excels in dynamic environments by offering real-time  
1429 knowledge updates and effective utilization of external knowledge sources with high interpretability.  
1430 However, it comes with higher latency and ethical considerations regarding data retrieval. On the  
1431 other hand, FT is more static, requiring retraining for updates but enabling deep customization of  
1432 the model's behavior and style. It demands significant computational resources for dataset preparation  
1433 and training, and while it can reduce hallucinations, it may face challenges with unfamiliar data.  
1434 In multiple evaluations of their performance on various knowledge-intensive tasks across different  
1435 topics, it was revealed that while unsupervised fine-tuning shows some improvement, RAG consistently  
1436 outperforms it, for both existing knowledge encountered during training and entirely new knowledge.  
1437 Additionally, it was found that LLMs struggle to learn new factual information through unsupervised  
1438 fine-tuning. The choice between RAG and FT depends on the specific needs for data dynamics,  
1439 customization, and computational capabilities in the application context. RAG and FT are not mutually  
1440 exclusive and can complement each other, enhancing a model's capabilities at different levels. In  
1441 some instances, their combined use may lead to optimal performance. The optimization process  
1442 involving RAG and FT may require multiple iterations to achieve satisfactory results.

#### 1443 5. RLHF (DPO)

1444 Reinforcement learning from human feedback (RLHF) is a variant of reinforcement learning (RL)  
1445 that learns from human feedback instead of relying on an engineered reward function. Building on  
1446 prior work on the related setting of preference-based reinforcement learning (PbRL), it stands at  
1447 the intersection of artificial intelligence and human-computer interaction. This positioning offers  
1448 a promising avenue to enhance the performance and adaptability of intelligent systems while also  
1449 improving the alignment of their objectives with human values. The training of large language  
1450 models (LLMs) has impressively demonstrated this potential in recent years, where RLHF played a  
1451 decisive role in directing the model's capabilities toward human objectives. In the RLHF setting, the  
1452 learning agent needs to solve an RL task without having access to a reward function. To this end,  
1453 the agent usually simultaneously learns an approximation of the reward function (via the assumed  
1454 utility function) and an RL policy. Therefore, a generic RLHF algorithm consists of repeating two  
1455 phases: (1) reward learning and (2) RL training. The first phase can itself be decomposed into two  
1456 main steps: (i) generate queries to ask the oracle, (ii) train a reward function approximator with the  
1457 answers provided by the oracle. The RL training part is more conventional and is usually directly  
based on running a deep RL algorithm using the currently trained reward function approximator.

1458 After learning a reward model, or, more commonly, interleaved with reward model learning, the next  
1459 step is to train a policy that maximizes the expected accumulated reward. This section will discuss  
1460 algorithms for policy learning, which can be categorized into two main techniques: adaptation of  
1461 conventional RL algorithms and direct policy optimization (DPO). Recent studies have also shown  
1462 that DPO is particularly sensitive to distribution shifts between the base model outputs and the  
1463 preference data. This sensitivity can lead to poor performance when there's a mismatch between  
1464 the training data of the base model and the preference dataset. Iterative DPO has been proposed to  
1465 address this issue. New responses are generated with the latest policy model, and a critique (can be  
1466 either a separate reward model or the same policy network in a self-rewarding setting) is used for  
1467 preference labeling in each iteration. This approach can help mitigate the distribution shift problem  
1468 and potentially improve DPO's performance. Lastly, the tests in the DPO paper were primarily  
1469 conducted on simple cases, including the IMDB dataset for controlled sentiment generation and the  
1470 Reddit dataset for summarization. The DPO loss function aimed to maximize the disparity between  
1471 desired and undesired responses. However, this approach could be problematic. It might lead to  
1472 simultaneous increases or decreases in the rewards for both desired and undesired responses, as long  
1473 as the difference between them grows.

## 1474 6. AGENTS

1475 Artificial Intelligence is entering a pivotal era with the emergence of LLM agents, intelligent enti-  
1476 ties powered by large language models (LLMs) capable of perceiving environments, reasoning about  
1477 goals, and executing actions. Unlike traditional AI systems that merely respond to user inputs, mod-  
1478 ern LLM agents actively engage with their environments through continuous learning, reasoning,  
1479 and adaptation. Compared to conventional agent systems, LLM-based agents have achieved gen-  
1480 erality across multiple dimensions, including knowledge sources, generalization capabilities, and  
1481 interaction modalities. LLM agents have found applications across diverse fields, including health-  
1482 care, biomedicine, law, and education. LLM agents introduce transformative solutions, such as  
1483 intelligent online tutoring systems, revolutionizing education accessibility. Planning capabilities are  
1484 a critical aspect of LLM agents' abilities, enabling them to navigate through complex tasks and  
1485 problem-solving scenarios with high accuracy. Effective planning is essential for deploying LLM  
1486 agents in real-world applications, where they must handle a diverse range of complex tasks and  
1487 scenarios. The planning capability of an LLM agent can be viewed from two perspectives: task  
1488 decomposition and feedback-driven iteration.

1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511