## PHYX: DOES YOUR MODEL HAVE THE "WITS" FOR PHYSICAL REASONING?

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

033

034

037

038

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

### **ABSTRACT**

Existing benchmarks fail to capture a crucial aspect of intelligence: physical reasoning, the integrated ability to combine domain knowledge, symbolic reasoning, and understanding of real-world constraints. To address this gap, we introduce PHYX: the first large-scale benchmark designed to assess models' capacity for physicsgrounded reasoning in visual scenarios. PHYX includes 3K meticulously curated multimodal questions spanning 6 reasoning types across 25 sub-domains and 6 core physics domains: thermodynamics, electromagnetism, mechanics, modern physics, optics, and wave & acoustics. In our comprehensive evaluation, even stateof-the-art models struggle significantly with physical reasoning. **GPT-04-mini**, Gemini-2.5-Pro, and GPT-5 achieve only 45.8%, 62.4%, and 65.2% accuracy respectively—performance gaps exceeding 10% compared to human experts. Our analysis exposes critical limitations in current models: over-reliance on memorized disciplinary knowledge, excessive dependence on mathematical formulations, and surface-level visual pattern matching rather than genuine physical understanding. We provide in-depth analysis through fine-grained statistics, detailed case studies, and multiple evaluation paradigms to thoroughly examine physical reasoning capabilities. To ensure reproducibility, we implement an evaluation protocol based on widely-used toolkits such as VLMEvalKit and lmms-eval, enabling one-click evaluation. All source code and data are available on our anonymous repository: anonymous.4open.science.

### 1 Introduction

Physics is the most fundamental and all-inclusive of the sciences.

- Richard Feynman

State-of-the-art models (Guo et al., 2025; OpenAI, 2024b; Team, 2025) can now basically solve Olympiad-level mathematical problems with human-competitive accuracy on benchmarks including AIME (MAA, 2024), GPQA (Rein et al., 2024), MATH-500 (Hendrycks et al., 2021), Olympiad-Bench (He et al., 2024), etc. Emerging multimodal large language models (MLLMs) like GPT-40 (OpenAI, 2024a) further offer promising pathways by combining visual understanding with reasoning capabilities. Recent advances in multimodal foundation models have spurred the development of benchmarks assessing disciplinary knowledge (Yue et al., 2024) and mathematical problems (Wang et al., 2024a; Zhang et al., 2024; Lu et al.). However, these evaluations overlook a critical dimension of machine intelligence: physical reasoning, the ability to integrate disciplinary knowledge, symbolic operations, and understanding of real-world constraints.

To address these gaps, we present **PHYX**, the first large-scale benchmark designed for evaluating physics-based reasoning via multimodal problem-solving with three core innovations: (1) 3,000 newly collected questions with realistic physical scenarios requiring integrated visual analysis and causal reasoning, (2) Expert-validated data design covering six fundamental physics domains with representative examples illustrated in Figure 2, and six distinct physical reasoning types, (3) Strict unified three-step evaluation protocols account for varying instruction-following capabilities across models and enable accurate assessment of reasoning. Each scenario undergoes rigorous validation by physics Ph.D. students to ensure scientific accuracy while eliminating dataset bias.

In addition to MLLMs, our benchmark supports evaluating LLMs by translating the images into text descriptions, thereby enabling an assessment of LLMs on these visually-grounded tasks. Our evalua-

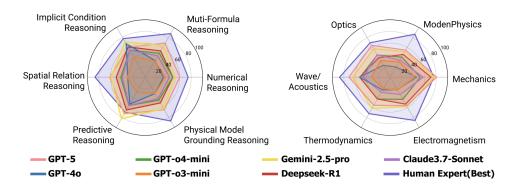


Figure 1: Accuracies of three leading MLLMs, two leading LLM and human performance on our proposed PHYX across 6 physical reasoning types and 6 domains.

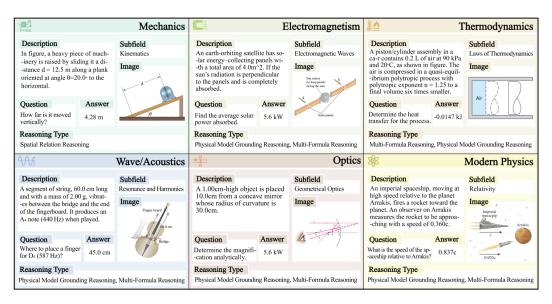


Figure 2: Sampled PHYX examples from each domain.

tion of 16 foundation models reveals an unprecedented capability gap: While the worst-performance group of physics undergraduates and graduates achieves 75.6% accuracy, the best-performing MLLM (GPT-5) scores only 65.2%. This 10-point performance chasm persists across all physics domains, most notably in Modern Physics (human 86.7% vs. model 56.5%) and Wave/Acoustics (human 86.7% vs. model 71.0%), as shown in Figure 1.

These results expose three critical shortcomings in current multimodal reasoning frameworks: (1) Visual reasoning errors (39.6%) indicate that models frequently misinterpret visual context, underscoring their limited capability in extracting and reasoning from physical scenarios. (2) The inconsistent performance across input variations: Full-Text, Text-DeRedundancy, and Text-Minimal, demonstrates that MLLMs remain overly dependent on textual descriptions, failing to effectively leverage visual input for reasoning. (3) Comparing physical reasoning performance to mathematical reasoning benchmarks such as MathVerse (Lu et al.) and MATH-V (Wang et al., 2024a) reveals that physical reasoning poses significantly greater challenges, highlighting a critical need for improved integration of abstract concepts and real-world knowledge. **PHYX** thus provides both a toolkit for model improvement and a roadmap for developing physically-grounded AI systems.

Our contributions can be summarized as follows: **Novel Benchmark Design:** We introduce PHYX, the first large-scale benchmark for evaluating the reasoning capabilities in the physical world for both multi-modal models and language models. Curated by experts, it spans 25 fine-grained domains and 6 reasoning types with realistic scenarios. **Versatile Evaluation Framework:** PHYX supports versatile evaluation frameworks, including *assessment formats* (multiple-choice vs. open-ended)

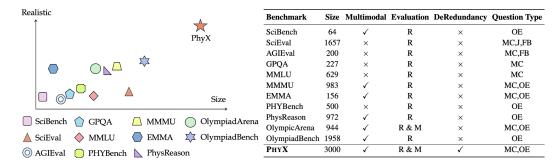


Figure 3: Comparison with existing benchmarks. Realistic refers to the extent to which the dataset contains visually realistic physical scenarios. Size indicates the number of physics questions with images in multimodal benchmarks or total physics questions in text-only benchmarks. For evaluation methods, R: rule-based, M: model-based. For question type, OE: Open-ended, MC: Multiple-choice, FB: Fill-in-the-blank, J: Judgement. In comparison, PHYX leads in all aspects.

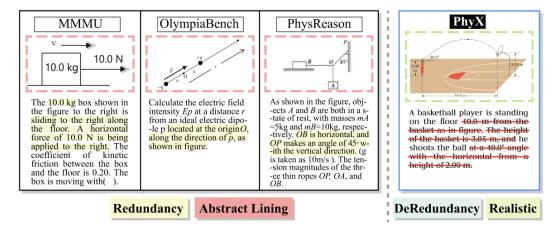


Figure 4: Existing benchmarks that contain physics questions suffer from information redundancy and abstract representation. In contrast, de-redundancy in PHYX increases the difficulty, as models can perceive concepts from ONE modality only. Additionally, realistic visuals challenges models to accurately apply physical laws.

and hierarchical answer judge (rule-based and model-based). It also seamlessly integrates with mainstream toolkits (e.g., VLMEvalKit, lmms-eval) for reproducible benchmarking. **Critical Insights on Reasoning:** We provide granular performance analysis and reveal some interesting observations, which sheds light on the design of the future models that jointly consider the disciplinary knowledge, symbolic operations, and real-world constraints for physical reasoning.

### 2 THE PHYX BENCHMARK

### 2.1 OVERVIEW OF PHYX

We introduce PHYX, a novel benchmark meticulously curated to assess the physical reasoning capabilities of foundation models. PHYX consists of 3,000 visually-grounded physics questions, meticulously curated to cover six distinct physics domains including *Mechanics* (550), *Electromagnetism* (550), *Thermodynamics* (500), *Wave/Acoustics* (500), *Optics* (500), and *Modern Physics* (400). Each problem in PHYX is centered around realistic physical scenarios to robustly assess the model's ability to reason the physical world. Detailed data statistics are summarized in Table 1, with representative question examples from each domains illustrated in Figure 2. To enable a comprehensive assessment, each question within PHYX has been categorized into six well-defined physical reasoning types: *Physical Model Grounding Reasoning*, *Spatial Relation Reasoning*, *Multi-Formula Reasoning*, *Implicit Condition Reasoning*, *Numerical Reasoning*, and *Predictive Reasoning*. Detailed definitions and illustrative examples of these reasoning types are provided in Appendix F.4.

Table 1: Key Statistics of PHYX.

Statistic	Number
Total new questions	6,000
- Multiple-choice questions	3,000 (50.0%)
- Open-ended questions	3,000 (50.0%)
Unique number of images	3,000
Unique number of questions	3,000
Maximum description length	288
Maximum question length	119
Maximum option length	46
Average description length	48.3
Average question length	14.6
Average option length	11.2

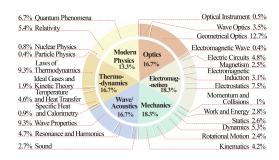


Figure 5: Fine-grained Distribution of PHYX.

Through its carefully curated structure and extensive coverage of diverse reasoning dimensions, PHYX represents a robust resource for systematically benchmarking and advancing the capabilities of foundation models in realistic physical reasoning tasks.

### 2.2 DATA CURATION PROCESS

**Data Collection.** To ensure high-quality data, we design a four-stage data collection process. Firstly, we conducted an in-depth survey of core physics disciplines to determine the coverage of our benchmark. We selected diverse physics domains and subfields, and defined a set of reasoning types. Secondly, we recruited a team of graduate students in STEM fields to serve as expert annotators. Annotators are instructed to comply with copyright and licensing rules by avoiding content from sources that restrict copying or redistribution. To mitigate potential data contamination in foundation models, they are advised to select questions for which answers are not immediately available alongside the problem, such as those found in separate materials or at the end of textbooks. Then, each openended question is required to be converted into a multiple-choice version, and vice versa. We also constructed three parallel versions of each question: (1) the original version; (2) a concise version where redundant textual information—those duplicated by the corresponding image—was removed; and (3) a question-only version that retains only the core question. Lastly, to support evaluation of LLMs, we used GPT-40 to generate descriptive captions for each image, aim to summarize the visual content in a self-contained textual form. This data curation process results in a diverse collection of 3,300 questions from various sources. The detailed annotation protocol is in Appendix I.

**Data Quality Control.** To further control the quality of our data, we perform a three-stage data cleaning process. First, we detect potentially duplicated questions by analyzing lexical overlap, followed by manual review from physics Ph.D. students to confirm and remove duplicates. Then, we filter out the shortest 10% of questions based on their textual length. This rigorous process plays a crucial role in maintaining the quality and difficulty of PHYX.

### 2.3 KEY DIFFERENCE COMPARED TO EXISTING BENCHMARKS

Compared with Scientific Knowledge Benchmarks. From Figure 3, science benchmarks like MMMU (Yue et al., 2024) cover broad disciplinary QA but lack focus on physical reasoning. These benchmarks often rely on basic understanding of disciplinary knowledge, with tasks that prioritize simple factual recall. In contrast, PHYX demands integration of visual cues with implicit physical laws, requiring models to perform context-driven inference. This targeted design evaluates multimodal reasoning about the physical world, exposing gaps in models' ability to handle scientific challenges.

Compared with Mathematical Reasoning Benchmarks. Mathematical reasoning benchmarks, such as MathVista (Lu et al.), MathVerse (Zhang et al., 2024), and MATH-V (Wang et al., 2024a), focus on logical deduction with clear expressions and explicit conditions, representing a subset of the challenges in physical reasoning. Physical reasoning extends beyond these by requiring models to model real-world contexts, identify implicit conditions from visual cues (e.g., Figure 6), and integrate the application of physical laws with symbolic logic, which are key capabilities absent in purely mathematical tasks. This makes PHYX a more comprehensive test of multimodal reasoning, capturing the complexity of real-world physics problems.

Table 2: Accuracy scores on the *testmini* subset of PHYX. The highest and the second highest scores of models in each setting are respectively highlighted in blue and red.

Models	Full	-Text	Text-DeR	edundancy	Text-Minimal	
Widels	Open-Ended	Multi-Choice	Open-Ended	Multi-Choice	Open-Ended	Multi-Choice
Random Choice	-	25	-	25	-	25
Human Expert (Worst)	-	-	75.6	-	-	-
Human Expert (Medium)	-	-	77.8	-	-	-
Human Expert (Best)	-	-	78.9	-	-	-
	M	ultimodal Large	Language Mode	els		
Claude3.7-Sonnet	44.4	65.8	42.2	64.5	17.2	41.6
Claude3.5-Sonnet	40.2	62.6	39.0	63.5	17.0	43.5
Claude3.5-Haiku	7.9	37.0	13.6	37.5	5.5	31.7
GPT-5	66.4	90.9	65.2	88.3	29.6	64.1
GPT-o4-mini	49.0	87.9	45.8	86.9	24.1	62.6
GPT-4o	33.9	61.0	32.5	57.6	14.3	43.8
Gemini-2.5-Pro	65.0	74.1	62.4	74.1	28.4	54.0
InternVL3-78B	35.9	45.6	33.1	46.9	14.8	40.5
QVQ-72B-Preview	17.5	40.0	17.2	40.9	7.6	33.1
Yi-VL-34B	3.5	34.8	3.4	34.1	1.9	34.1
InternVL3-14B	9.0	46.9	7.9	47.5	5.1	45.9
InternVL3-8B	6.3	45.5	6.5	44.9	4.6	44.0
MiniCPM-o-8B	7.1	31.8	7.2	31.6	3.2	34.2
LLaVA-OneVision-7B	7.2	37.7	5.7	37.3	2.7	38.0
DeepSeek-VL2-4.5B	11.4	28.2	10.2	27.8	4.7	27.3
Kimi-VL-A3B-Instruct-2.8B	15.6	37.1	15.4	38.7	8.1	39.3
Kimi-VL-A3B-Thinking-2.8B	15.3	34.4	15.8	33.2	7.4	27.0
		Large Langu	age Models			
DeepSeek-R1	51.8	63.1	51.2	62.9	22.2	43.6
DeepSeek-V3	40.7	70.8	36.3	67.5	16.2	49.9
GPT-o3-mini	36.9	78.5	31.5	76.9	14.3	56.2
Qwen3-4B	29.6	49.2	27.5	48.4	12.1	41.8
Qwen3-8B	32.2	50.4	31.6	48.8	13.0	37.2
Owen3-14B	35.3	57.2	33.3	56.5	13.6	44.4

Compared with Physics-related Benchmarks Existing benchmarks (e.g., PHYBench (Qiu et al., 2025), UGPhysics (Xu et al., 2025a), OlympiadBench (He et al., 2024)) prioritize text-based problems or schematic visuals, limiting their assessment of multimodal reasoning. In detail, PHYBench's problems and UGPhysics's questions rely heavily on textual descriptions, while OlympiadBench's problems use simplified diagrams, as shown in Figure 4. These benchmarks mainly test disciplinary knowledge but overlook the integration of visual perception with implicit physical constraints. PHYX bridges these gaps by embedding high-fidelity visual scenarios that require models to decode complex visual cues, infer context-specific physical laws and then reason about problems. Additionally, PHYX mandates equal reliance on both modalities with information de-redundancy, providing a rigorous evaluation of professional-level physical reasoning in MLLMs.

### 3 EXPERIMENTS

### 3.1 EXPERIMENTAL SETUP

**The testmini Subset.** PHYX comprises 3,000 high-quality visual physics problems and 18,000 corresponding test instances. To streamline evaluation, we extract a smaller representative subset named *testmini* including 1,000 problems and 6,000 instances. The construction of *testmini* involved a proportional random sampling strategy across different physics domains of PHYX. The quantitative evaluations in all subsequent experiments were conducted on this *testmini* subset, while the full set of results on all 3,000 problems are provided in Appendix E.

**Baselines.** We include random choice as a naive baseline, and we recruit 15 undergraduate and graduate physics students to represent the expert performance baseline, each student was tasked with completing 18 questions. The students were divided into three groups of five, and the results of each group are reported separately. Then, we conduct experiments on (a) Reasoning MLLMs: Gemini-2.5-Pro (Team, 2025), GPT-5 (OpenAI, 2025c), GPT-04-mini (OpenAI, 2025b), Claude-3.7-Sonnet (claude, 2025), QVQ-Preview (Team, 2024), LLaVA-OneVision (Li et al., 2024), MiniCPM-

Table 3: Average scores by model across different domains of physics with *Open-Ended Text-DeRedundancy* questions. The highest and the second highest scores of models in each domain are respectively highlighted in each setting in red and blue.

Models	Overall	Mechanics	Electro- magnetism	Thermo- dynamics	Waves & Acoustics	Optics	Modern Physics
Human Expert (Worst)	75.6	76.5	60.0	66.7	86.7	69.2	86.7
Human Expert (Medium)	77.8	94.1	53.3	60.0	93.3	76.9	86.7
Human Expert (Best)	78.9	76.5	86.7	73.3	86.7	69.2	86.7
		Multimodal I	Large Language	Models			
Claude3.7-Sonnet	42.2	58.2	36.7	31.5	46.7	44.6	35.2
Claude3.5-Sonnet	39.0	53.5	27.8	33.3	49.7	35.5	3.9
Claude3.5-Haiku	13.6	18.8	8.9	11.5	18.8	12.0	11.5
GPT-5	65.2	80.2	57.4	55.2	71.0	63.3	56.5
GPT-o4-mini	45.8	52.3	43.2	41.8	52.7	44.0	40.6
GPT-4o	32.5	45.9	24.3	26.1	53.9	23.5	21.2
Gemini-2.5-Pro	62.4	77.6	59.2	61.8	64.8	57.2	53.3
InternVL3-78B	33.1	48.8	27.2	25.5	43.0	28.9	24.8
QVQ-72B-Preview	17.2	31.7	11.2	10.9	20.0	12.0	16.9
Yi-VL-34B	3.4	1.8	3.5	4.8	2.4	4.2	3.6
InternVL3-14B	7.9	12.4	8.9	4.2	8.5	4.8	8.5
InternVL3-8B	6.5	10.6	6.5	3.6	4.9	6.6	6.7
MiniCPM-o-8B	7.2	11.8	6.5	6.1	7.3	6.0	5.5
LLaVA-OneVision-7B	5.7	10.6	4.1	6.1	7.3	3.0	3.0
DeepSeek-VL2-4.5B	10.2	16.5	7.1	10.3	13.3	9.0	4.8
Kimi-VL-A3B-Instruct-2.8B	15.4	20.6	10.1	13.3	20.0	16.2	12.1
Kimi-VL-A3B-Thinking-2.8B	15.8	25.9	15.4	7.9	20.6	13.3	11.5
Large Language Models							
DeepSeek-R1	51.2	71.8	53.2	41.8	53.9	39.8	46.1
DeepSeek-V3	36.3	52.9	39.6	28.5	36.4	28.9	30.9
GPT-o3-mini	31.5	41.8	24.9	23.6	32.1	33.7	32.7
Qwen3-4B	27.5	42.9	23.7	21.2	35.8	21.1	20.0
Qwen3-8B	31.6	51.2	26.6	19.4	37.0	29.5	25.5
Qwen3-14B	33.3	52.9	30.8	18.2	40.0	27.1	30.3

o (Yao et al., 2024), Kimi-VL-A3B-Thinking (Team et al., 2025), (b) General MLLMs: GPT-4o (OpenAI, 2024a), Claude-3.5-Sonnet (claude, 2024b), Claude-3.5-Haiku (claude, 2024a), InternVL3 (Zhu et al., 2025), Yi-VL (Young et al., 2024), Kimi-VL-A3B-Instruct (Team et al., 2025), (c) LLMs: GPT-o3-mini (OpenAI, 2025a), DeepSeek-R1 (Guo et al., 2025), DeepSeek-V3 (DeepSeek-AI, 2025), Qwen3 (Yang et al., 2025), augmented with image captions generated by GPT-4o.

### 3.2 EVALUATION PROTOCOLS

Our evaluation is conducted with Chain-of-Thought (CoT) prompting to assess the reasoning capability of models. For both open-ended (OE) and multiple-choice (MC) questions, the instruction-following capabilities of models can vary significantly. To this end, we design a universal evaluation pipeline for all recent LLMs and MLLMs with different instruction-following capabilities:

**Step 1. Prediction Generation.** Initially, the models generate predictions given the input query, which incorporates different problem descriptions according to the specific settings, the question, and the image, using the template defined in Appendix G.1.

**Step 2. Answer Extraction.** The raw predictions often contain reasoning steps, explanations, or irrelevant conversational filler. To precisely extract the definitive answer from these raw outputs, we separately employ rule-based answer extraction strategies, which are detailed in Appendix G.2.

**Step 3. Answer Judgment.** For OE questions, the next step is comparing the extracted answer against the ground truth. Given that answers in OE physics questions can be expressed in myriad ways, we propose an evaluation mechanism using an LLM, such as DeepSeek-V3 (DeepSeek-AI, 2025), as a judge, using the template defined in Appendix G.3. We feeds the answer extracted and the ground truth to an LLM multiple times and checks if an LLM succeed in all attempts. A preliminary study of 200 examples shows that DeepSeek-V3 can judge the answer with more than 99% accuracy with an affordable cost. For MC questions, we first attempt to directly match the option letter. If it fails, we then use an LLM as a judge, using the template for OE questions.

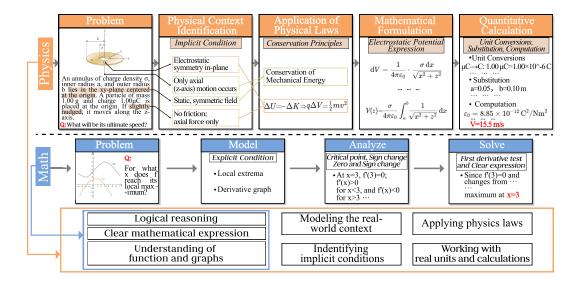


Figure 6: An real example of reasoning trajectory based on GPT-40 and the comparison of required capabilities when solving physical and mathematical problems.

### 3.3 MAIN RESULTS

In this section, we present a comprehensive comparison of LLMs and MLLMs on the PHYX benchmark, detailed in Table 2 and Table 3. Our key findings can be summarized as follows:

**Challenging Nature of PHYX.** PHYX presents significant challenges for current models. Notably, even the worst human experts achieve accuracy of 75.6%, significantly outperforming all the models included in our comparative analysis. This disparity demonstrates an existing gap between human expertise and current model capabilities, reflecting the demanding standards inherent in PHYX.

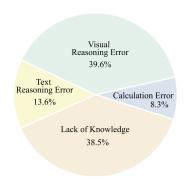
**Question Format as a Determinant of Model Discriminability.** The result reveals that multiple-choice questions reduce the performance gap across models, enabling weaker models to rely on surface-level cues. In contrast, open-ended questions demand genuine reasoning and precise answer generation, leading to greater differentiation between models. This suggests that the open-ended format provides higher discriminative power when evaluating multimodal reasoning capabilities.

**Domain-specific Variability in Multimodal Reasoning.** As shown in Table 3, in domains such as Waves/Acoustics and Mechanics, which typically include natural images and questions requiring relatively less reasoning, models tend to achieve higher performance. Conversely, in domains such as Thermodynamics and Modern Physics, where tasks frequently demand intricate visual perception and multi-step reasoning, models performance is generally lower.

**Reasoning-oriented Models Perform Better.** Leading reasoning-oriented models such as GPT-o4-mini and DeepSeek-R1 achieve accuracies of 45.8% and 51.2%, respectively, significantly outperforming general-purpose models like GPT-40 and Claude3.7-Sonnet. The results highlight the advantage of models specifically optimized for reasoning tasks, suggesting that architectural and training differences play a key role in bridging the multimodal reasoning gap.

### 3.4 DISCUSSION

Mathematical Reasoning Is Not Enough. Comparing GPT-4o's performance on PHYX to its previously reported results on MathVista (63.8%) and MATH-V (63.8%), we observe substantially lower accuracy in physical reasoning tasks, underscoring that these tasks present challenges that go beyond mathematical abstraction. As shown in Figure 6, unlike mathematics problems, where symbolic manipulation and abstraction are often sufficient, physical reasoning requires models to decode implicit conditions in the problem statement (e.g., interpreting "smooth surface" as implying zero friction), ground physical laws in concrete visual and material contexts (e.g., recognizing whether a wooden block will float or whether a surface is rough or smooth), and maintain internal



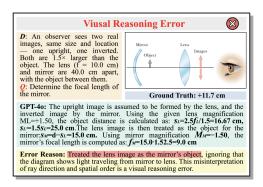


Figure 7: The error distribution over 90 annotated errors with a typical visual reasoning error, which is easy for humans but challenging for GPT-40. More examples can be found in the Appendix H.

consistency across reasoning chains, since the laws of physics remain invariant regardless of the reasoning trajectory. This tight integration of perception, abstraction, and dynamic simulation makes physical reasoning qualitatively more demanding, revealing a fundamental barrier that current LLMs and MLLMs have yet to overcome.

Impact of Redundancy Removal on Models. We observe that most MLLMs exhibit performance degradation under the Text-DeRedundancy setting, reflecting that most MLLMs are built upon a powerful language model backbone, and thus tend to over-rely on textual information while under-utilizing visual input. However, certain models (e.g., Claude 3.5-Haiku) benefit from redundancy removal, which aligns with recent findings (Li et al., 2025; Mao et al., 2025) that imperfect cross-modal alignment can introduce noisy or conflicting signals. In such cases, redundant textual cues exacerbate alignment issues, and their removal clarifies multimodal inputs. Notably, this effect is unique to MLLMs and does not occur in pure-text LLMs, further suggesting that the root cause lies in cross-modal fusion rather than general reasoning ability.

MLLMs' Physical Reasoning Relies More on Text. Our experiments show a clear performance gradient across the three input variations: Full-Text, Text-DeRedundancy, and Text-Minimal, with decreasing accuracy in that order. This indicates that most MLLMs rely heavily on detailed textual descriptions, highlighting their limited ability to reason purely from visual context.

Competitive LLMs Highlight Limitations in Multimodal Fusion. Despite lacking direct visual input, LLMs such as DeepSeek-R1 and GPT-o3-mini perform competitively with most multimodal models. The strong performance of LLMs suggests that, in many cases, the caption provides sufficient visual context for reasoning. This highlights both the impressive generalization capabilities of LLMs and the current limitations of MLLMs in leveraging raw visual signals for physical reasoning.

### 3.5 ERROR ANALYSIS

To dive into the reasoning capabilities and limitations of models, we meticulously inspected 96 randomly sampled incorrect predictions and performed an in-depth analysis based on GPT-40. The objectives of this analysis were twofold: to identify current model weaknesses and to guide future enhancements in model design and training. The distribution of these errors is illustrated in Figure 7, and a comprehensive case study of 30 notable cases is included in Appendix H.

**Reasoning Errors** (53.2%) encompass both visual and textual reasoning failures. While *visual reasoning errors* typically arise from incorrect extraction, spatial relationships, or reasoning based on visual information (e.g., misreading the voltage in Appendix 8), *textual reasoning errors* are characterized by misinterpretation of textual content, such as overlooking explicit conditions (e.g., ignoring the instruction to neglect friction in Appendix 4). Furthermore, our analysis reveals three deeper challenges in multimodal reasoning:

- *Context Switching*. Rapidly transitioning between textual and visual modalities can cause models to lose focus or misinterpret key data. Prior work (Zhang et al., 2025; Li et al., 2025) confirms that modality switching introduces significant cognitive load, leading to unstable attention, misalignment in cross-modal representations, and disrupted reasoning chains, especially in tasks requiring deep integration of visual and textual information.

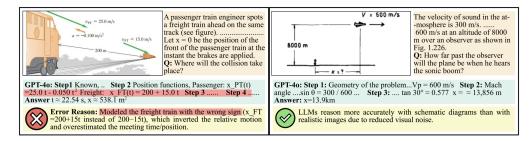


Figure 8: Both are relative motion problems, easy for human experts. *Left*: Error case from PHYX with realistic images, which introduce visual noise and make reasoning harder. *Right*: Correct case from concurrent work SeePhys (Xiang et al., 2025), where schematics image simplify perception.

- *Cross-Modal Reasoning Difficulty.* Models often struggle to deeply integrate visual and textual modalities, particularly when fine-grained visual cues must align with detailed textual descriptions. Existing studies (Mao et al., 2025; Yue et al., 2024) highlight that current fusion strategies frequently yield superficial alignment, limiting cross-modal reasoning.
- *Visual Realism Challenge*. As illustrated in Figure 8, unlike other benchmarks that primarily employ abstract line drawings or schematic sketches, our dataset contains realistic images. This realism significantly increases perceptual and reasoning difficulty: models must construct an abstract physical problem from a noisy, unstructured visual input rather than relying on simplified shapes.

**Lack of Knowledge (38.5%)** reflects GPT-4o's incomplete understanding of physical domain knowledge. As shown in Appendix 25, the model lacks the fundamental knowledge. Specifically, it ignores that the slower speed in the liver requires a correction when estimating depth from the reflection geometry, leading to an overestimated result.

**Calculation Error (8.3%)** refer to mistakes in arithmetic operations, formula application, or unit conversions. These errors indicate that the model has grasped the physical context and relevant concepts but fails in the final step of numerical computation.

### 4 RELATED WORK

Multi-modal Large Language Models. Multi-modal large language models (MLLMs) (OpenAI, 2025b; Team, 2025) have shown great potential and achieved excellent visual understanding by integrating visual and textual data across a wide range of multimodal tasks. Recent advances in LLMs have motivated efforts (Wei et al., 2022; Ouyang et al., 2022) to explore MLLM reasoning. Despite such achievements, it remains unclear whether these models truly possess advanced reasoning abilities, especially in the physical area that is closer to the real world. To bridge this gap and evaluate the physical reasoning capabilities of MLLMs, we introduce PHYX, a multimodal benchmark to evaluate the real reasoning ability of recent advanced MLLMs in physics.

MLLM Benchmarks. Recently, several MLLM scientific benchmark (Yue et al., 2024; Wang et al., 2024b; He et al., 2024; Huang et al., 2024; Zhang et al., 2025; Hao et al., 2025) have also been proposed. For example, PhysReason (Zhang et al., 2025) includes a multimodal subset of 972 physics problems with figures to evaluate the MLLMs. EMMA (Hao et al., 2025) comprises 2,788 problems covering various scientific areas such as mathematics, physics, and coding. However, all of these benchmarks contain only a small subset of data in physics, which still cannot fully evaluate MLLM's ability on physical reasoning. More related works are discussed in Appendix D.

### 5 CONCLUSION

Existing benchmarks have overlooked the critical task of physical reasoning, which requires integrating domain knowledge, symbolic reasoning, and real-world constraints. To address this, we present PHYX, the first large-scale benchmark for evaluating physical reasoning in multimodal, visually grounded scenarios. Through rigorous evaluation, we reveal that state-of-the-art models exhibit significant limitations in physical reasoning. Our findings highlight the urgent need for future models to improve deep physical reasoning over surface-level associations.

### REFERENCES

486

487

488

489

490 491

492

495

496 497

498

499

500

501

502

504

505

506

507

508

509

510

511

512

513

514515

516

517

518

519

520

521

522

523

524

525526

527

528 529

530

531

532

534

535

536

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- claude. Claude 3.5 haiku. https://www.anthropic.com/claude/haiku, 2024a.
- claude. Introducing claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet, 2024b.
  - claude. Claude 3.7 sonnet and claude code. https://www.anthropic.com/news/claude-3-7-sonnet, 2025.
  - DeepSeek-AI. Deepseek-v3 technical report, 2025. URL https://arxiv.org/abs/2412.19437.
    - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
    - Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025.
    - Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3828–3850, 2024.
    - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
    - Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv* preprint arXiv:2103.03874, 2021.
    - Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, et al. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent ai. *Advances in Neural Information Processing Systems*, 37:19209–19253, 2024.
    - Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.
    - Mingxiao Li, Na Su, Fang Qu, Zhizhou Zhong, Ziyang Chen, Yuan Li, Zhaopeng Tu, and Xiaolong Li. Vista: Enhancing vision-text alignment in mllms via cross-modal mutual information maximization. *arXiv* preprint arXiv:2505.10917, 2025.
    - Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*.
  - MAA. American invitational mathematics examination aime. In Invitational Mathematics Examination AIME2024, Febru-American 2024. URL https://maa.org/math-competitions/ ary american-invitational-mathematics-examination-aime.
    - Song Mao, Yang Chen, Pinglong Cai, Ding Wang, Guohang Yan, Zhi Yu, and Botian Shi. Investigating redundancy in multimodal large language models with multiple vision encoders. *arXiv* preprint *arXiv*:2507.03262, 2025.

- OpenAI. Gpt-4o system card, 2024a. URL https://arxiv.org/abs/2410.21276.
- OpenAI. Learning to reason with llms, 2024b. URL https://openai.com/index/learning-to-reason-with-llms/.
- OpenAI. Openai o3-mini: Pushing the frontier of cost-effective reasoning. https://openai.com/index/openai-o3-mini/, 2025a.
  - OpenAI. Introducing openai o3 and o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/, 2025b.
  - OpenAI. Gpt-5 system card. https://cdn.openai.com/gpt-5-system-card.pdf, August 13 2025c.
  - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
  - Shi Qiu, Shaoyang Guo, Zhuo-Yang Song, Yunbo Sun, Zeyu Cai, Jiashen Wei, Tianyu Luo, Yixuan Yin, Haoxu Zhang, Yi Hu, et al. Phybench: Holistic evaluation of physical perception and reasoning in large language models. *arXiv preprint arXiv:2504.16074*, 2025.
  - David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
  - Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. Scieval: A multi-level large language model evaluation benchmark for scientific research. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 19053–19061, 2024.
  - Gemini Team. Gemini 2.5: Our most intelligent ai model, 2025.

    URL https://blog.google/technology/google-deepmind/
    gemini-model-thinking-updates-march-2025/#gemini-2-5-thinking.
  - Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, Congcong Wang, Dehao Zhang, Dikang Du, Dongliang Wang, Enming Yuan, Enzhe Lu, Fang Li, Flood Sung, Guangda Wei, Guokun Lai, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haoning Wu, Haotian Yao, Haoyu Lu, Heng Wang, Hongcheng Gao, Huabin Zheng, Jiaming Li, Jianlin Su, Jianzhou Wang, Jiaqi Deng, Jiezhong Qiu, Jin Xie, Jinhong Wang, Jingyuan Liu, Junjie Yan, Kun Ouyang, Liang Chen, Lin Sui, Longhui Yu, Mengfan Dong, Mengnan Dong, Nuo Xu, Pengyu Cheng, Qizheng Gu, Runjie Zhou, Shaowei Liu, Sihan Cao, Tao Yu, Tianhui Song, Tongtong Bai, Wei Song, Weiran He, Weixiao Huang, Weixin Xu, Xiaokun Yuan, Xingcheng Yao, Xingzhe Wu, Xinxing Zu, Xinyu Zhou, Xinyuan Wang, Y. Charles, Yan Zhong, Yang Li, Yangyang Hu, Yanru Chen, Yejie Wang, Yibo Liu, Yibo Miao, Yidao Qin, Yimin Chen, Yiping Bao, Yiqin Wang, Yongsheng Kang, Yuanxin Liu, Yulun Du, Yuxin Wu, Yuzhi Wang, Yuzi Yan, Zaida Zhou, Zhaowei Li, Zhejun Jiang, Zheng Zhang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Zijia Zhao, and Ziwei Chen. Kimi-VL technical report, 2025. URL https://arxiv.org/abs/2504.07491.
  - Qwen Team. Qvq: To see the world with wisdom, December 2024. URL https://qwenlm.github.io/blog/qvq-72b-preview/.
  - Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024a.
  - Lintao Wang, Encheng Su, Jiaqi Liu, Pengze Li, Peng Xia, Jiabei Xiao, Wenlong Zhang, Xinnan Dai, Xi Chen, Yuan Meng, et al. Physunibench: An undergraduate-level physics reasoning benchmark for multimodal models. *arXiv preprint arXiv:2506.17667*, 2025.

- Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. In *International Conference on Machine Learning*, pp. 50622–50649. PMLR, 2024b.
- 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, 2022.
- Kun Xiang, Heng Li, Terry Jingchen Zhang, Yinya Huang, Zirong Liu, Peixin Qu, Jixi He, Jiaqi Chen, Yu-Jie Yuan, Jianhua Han, et al. Seephys: Does seeing help thinking?—benchmarking vision-based physics reasoning. *arXiv preprint arXiv:2505.19099*, 2025.
- Xin Xu, Qiyun Xu, Tong Xiao, Tianhao Chen, Yuchen Yan, Jiaxin Zhang, Shizhe Diao, Can Yang, and Yang Wang. Ugphysics: A comprehensive benchmark for undergraduate physics reasoning with large language models. *arXiv preprint arXiv:2502.00334*, 2025a.
- Yinggan Xu, Yue Liu, Zhiqiang Gao, Changnan Peng, and Di Luo. Physense: Principle-based physics reasoning benchmarking for large language models. *arXiv preprint arXiv:2505.24823*, 2025b.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pp. 169–186. Springer, 2024.
- Xinyu Zhang, Yuxuan Dong, Yanrui Wu, Jiaxing Huang, Chengyou Jia, Basura Fernando, Mike Zheng Shou, Lingling Zhang, and Jun Liu. Physreason: A comprehensive benchmark towards physics-based reasoning. *arXiv preprint arXiv:2502.12054*, 2025.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL https://arxiv.org/abs/2311.07911.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv* preprint arXiv:2504.10479, 2025.

### **Table of Contents in Appendix**

A The Use of Large Language Models **Reproducibility Statement Ethics Statement More Related Work** D **More Experimental Results More Dataset Details** F.2 F.3 **G** More Evaluation Details H Case Study **Data Annotation Protocol** I.1 I.2 I.3 I.4 I.5 I.6 I.7

### A THE USE OF LARGE LANGUAGE MODELS

In preparing this manuscript, we used a Large Language Model (LLM) solely to assist with minor language polishing and improvements in readability. The LLM did not contribute to research ideation, analysis, or substantive writing. All scientific content and conclusions are entirely the responsibility of the authors.

### B REPRODUCIBILITY STATEMENT

We have made every effort to ensure the reproducibility of our results. All code, scripts, and dataset used in this work are available in an anonymous repository anonymous.4open.science. Detailed descriptions of the benchmark setup, data collection, and preprocessing procedures are provided in Section 3 of the main text and in Appendix G. With these resources, all reported results can be fully reproduced.

### C ETHICS STATEMENT

**Legal Compliance.** All questions included in PHYX are sourced from publicly accessible materials. During data collection, annotators are instructed to strictly follow the copyright and licensing terms of the original platforms. Any content from sources that prohibit reuse or redistribution MUST be explicitly excluded. PHYX is a non-commercial project, and its usage aligns with the principles outlined in Fair Use §107: "the fair use of a copyrighted work, including such use by ...... scholarship, or research, is not an infringement of copyright", where fair use is determined by "the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes" and "the effect of the use upon the potential market for or value of the copyrighted work."

**Dataset Intended Usage and License.** The full details of the PHYX dataset are presented in this paper, and both the PHYX and code for reproducing results will be made publicly available. The PHYX dataset is not supposed to be used to train models for cheating. The primary goal is to support the research community in benchmarking and advancing physical reasoning in LLMs and MLLMs. We take full responsibility for any rights violation that may arise. Both the PHYX data and our open-source code are released under the MIT license.

### D MORE RELATED WORK

Several LLM benchmarks (Hendrycks et al.; Sun et al., 2024; Rein et al., 2024; Austin et al., 2021; Zhou et al., 2023) have been proposed to evaluate LLM's ability on various aspects. Among these works, the most related one is PHYBench (Qiu et al., 2025), which also focuses in the physic reasoning area. Although evaluating the same discipline, their scope remains narrow since it includes only a small number of questions, making it insufficient to fully assess a model's reasoning capabilities. Furthermore, PHYBench concentrates exclusively on evaluating through text. However, in real-world scenarios, solving physics problems also requires visual perception and interpretation. Concurrently, three related efforts (Wang et al., 2025; Xiang et al., 2025; Xu et al., 2025b) have emerged. While sharing a similar motivation, these benchmarks rely on schematic images that simplify visual perception. In contrast, our dataset leverages realistic images, which introduce visual noise and thereby make the reasoning process more challenging and closer to real-world conditions.

### E MORE EXPERIMENTAL RESULTS

Table 4 reports the accuracy scores of the leading MLLM (GPT-40) on the full *test* subset. The minor differences between overall scores on the *test* subset and the *testmini* subset, suggest that the *testmini* subset effectively mirrors the *test* subset, serving as a valuable evaluation subset for model development, especially for those who have limited computing resources.

Table 4: Accuracy scores by GPT-40 across different domains of physics with open-ended text de-redundancy questions on the *test* subset and *testmini* subset of PHYX.

GPT-40	Overall	Mechanics	Electro- magnetism	Thermo- dynamics	Waves & Acoustics	Optics	Modern Physics
testmini	32.5	45.9	24.3	26.1	53.9	23.5	21.2
test	34.0	45.8	37.1	20.0	38.2	27.2	34.3

### F MORE DATASET DETAILS

### F.1 QUESTION DISTRIBUTION

All questions in PHYX are written in English. Figure 9 presents the distribution of word counts of questions in the Text-DeRedundancy setting, demonstrating the variation in question lengths. The similarity between the median and average word counts suggests a roughly symmetrical distribution.

### F.2 Introduction of Domain and Subfield

As shown in Table 5, PHYX covers 6 core domains and 25 subdomains.

**Mechanics.** Mechanics is the branch of physics concerned with the motion of objects and the forces that cause or change this motion. It encompasses both classical mechanics and key subfields such as *Kinematics* (e.g., velocity, acceleration, free fall), *Dynamics* (e.g., Newton's laws, force analysis, friction), *Work and Energy* (e.g., work-energy theorem, mechanical energy conservation), *Momentum and Collisions* (e.g., conservation of momentum, elastic and inelastic collisions), *Rotational Motion* (e.g., torque, angular acceleration, moment of inertia), and *Statics* (e.g., torque balance, structural analysis). Mechanics lays the groundwork for much of physics, enabling the understanding of how and why objects move or remain at rest in various physical systems.

**Electromagnetism.** Electromagnetism explores the interactions between electric charges and magnetic fields. It includes the subfields of *Electrostatics* (e.g., Coulomb's law, electric fields and potential), *Electric Circuits* (e.g., Ohm's law, circuit analysis, RC circuits), *Magnetism* (e.g., magnetic fields, Lorentz force, Ampère's law), *Electromagnetic Induction* (e.g., Faraday's law, Lenz's law, inductance), and optionally *Maxwell's Equations and Electromagnetic Waves* for advanced topics. This domain underpins much of modern technology, including electric circuits, motors, and wireless transmission.

**Thermodynamics.** Thermodynamics is the study of heat, energy, and their transformations. Its subtopics include *Temperature and Heat Transfer* (e.g., conduction, convection, radiation), *Specific Heat and Calorimetry* (e.g., phase changes, heat calculations), *Laws of Thermodynamics* (e.g., energy conservation, entropy), and *Ideal Gases and Kinetic Theory* (e.g., gas laws, internal energy, pressure). This domain is central to engines, thermal systems, and understanding natural processes.

**Wave/Acoustics.** This domain investigates wave behavior and sound phenomena. Core subfields include *Wave Properties* (e.g., speed, frequency, wavelength, interference), *Sound* (e.g., pitch, loudness, Doppler effect, standing waves), and *Resonance and Harmonics* (e.g., resonant frequencies, vibrations in strings and air columns). These concepts are crucial in fields ranging from acoustics to telecommunications.

**Optics.** Optics studies the behavior and properties of light. It includes *Geometrical Optics* (e.g., reflection, refraction, lens imaging, total internal reflection), *Wave Optics* (e.g., interference, diffraction, polarization), and *Optical Instruments* (e.g., microscopes, telescopes, image formation). Optics has broad applications in imaging, vision science, and photonics.

**Modern Physics.** Modern Physics addresses phenomena beyond the scope of classical mechanics. Its key subfields include *Relativity* (e.g., time dilation, mass-energy equivalence), *Quantum Phenomena* (e.g., photoelectric effect, atomic models), *Nuclear Physics* (e.g., radioactivity, nuclear reactions, mass defect), and optionally *Particle Physics* (e.g., elementary particles, the Standard Model). These topics form the theoretical basis of contemporary physics and technology.

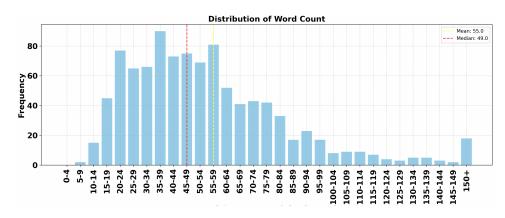


Figure 9: The distribution of the number of words per question in PHYX.

Domain	Subfields			
Optics	Optical Instrument, Wave Optics, and Geometrical Optics			
Electromagnetism	Electromagnetic Wave, Electric Circuits, Magnetism, Electromagnetic Induction, and Electrostatics			
Mechanics	Momentum and Collisions, Work and Energy, Statics, Dynamics, Relational Motion, and Kinematics.			
Wave/Acoustics	Sound, Resonance and Harmonics, and Wave Properties			
Thermodynamics	Specific Heat and Calorimetry, Temperature and Heat Transfer, Ideal Gases and Kinetic Theory, and Laws of Thermodynamics			
Modern Physics	Particle Physics, Nuclear Physics, Relativity, and Quantum Phenomena			

Table 5: Subfields included in each domain in PHYX.

### F.3 IMAGES BY DOMAINS

In this section, we present example images from the physics problems in PHYX. Figure 10, Figure 11, Figure 12, Figure 13, Figure 14 and Figure 15 show images from problems under the categories of Mechanics, Electromagnetism, Thermodynamics, Wave/Acoustics, Optics, and Modern Physics, respectively.

We observe that the images in our dataset are highly realistic, often depicting concrete physical scenarios rather than stylized or abstract illustrations. While they are not real-world photographs, these visuals are grounded in plausible physical settings. This realism provides essential context for physical reasoning and helps bridge the gap between abstract physics principles and their real-world manifestations.

Across domains, the visual characteristics vary in alignment with the nature of the physical concepts. Despite their domain-specific variations, a unifying theme across all categories is the consistent use of realistic and context-rich imagery, which provides essential grounding for physical interpretation and distinguishes our benchmark from other datasets with overly synthetic or schematic visual content.

### F.4 PHYSICAL REASONING DEFINITION

Six physical reasoning types are defined in Table 6.

### G MORE EVALUATION DETAILS

We conduct all experiments on NVIDIA A100 80G GPUs.

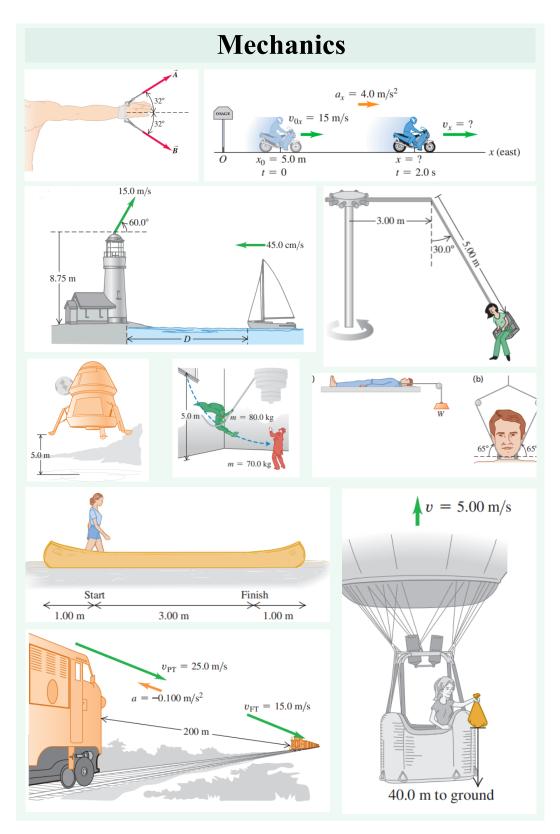


Figure 10: Examples of the visual context for the Mechanics domain.

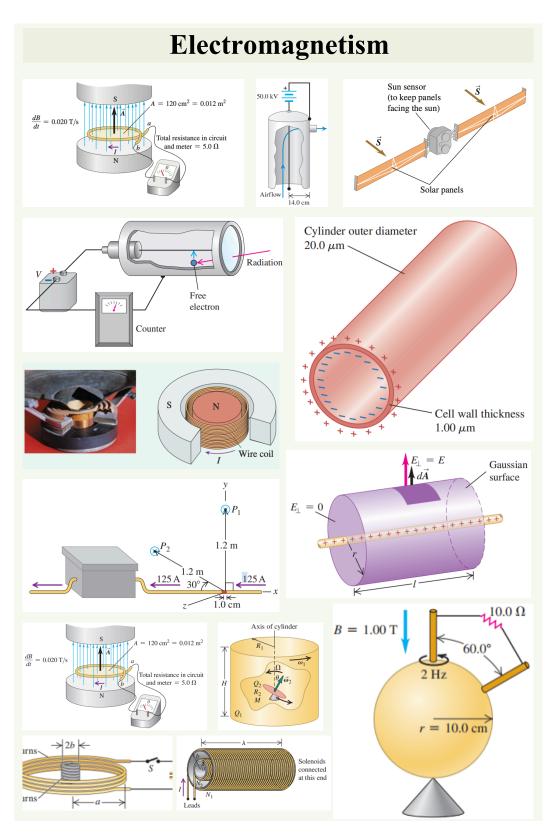


Figure 11: Examples of the visual context for the *Electromagnetism* domain.

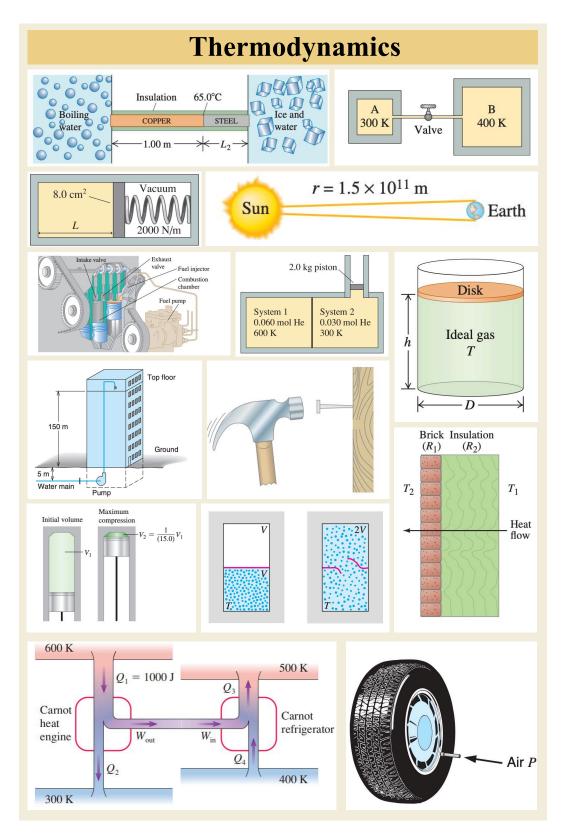


Figure 12: Examples of the visual context for the *Thermodynamics* domain.

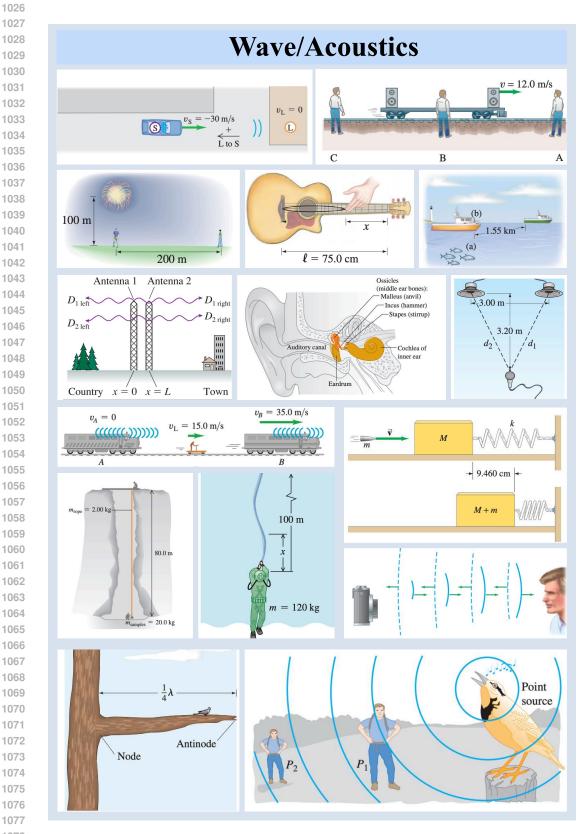


Figure 13: Examples of the visual context for the Wave/Acoustics domain.

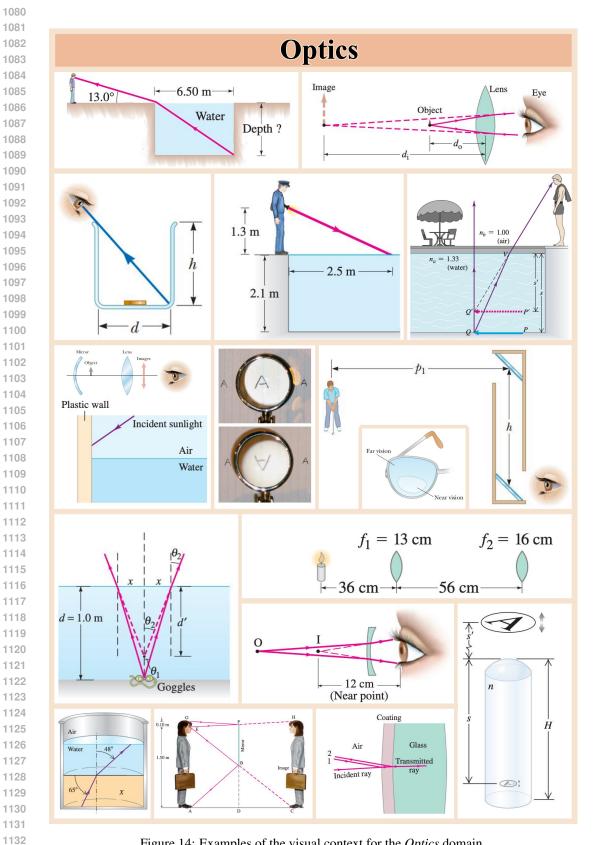


Figure 14: Examples of the visual context for the *Optics* domain.

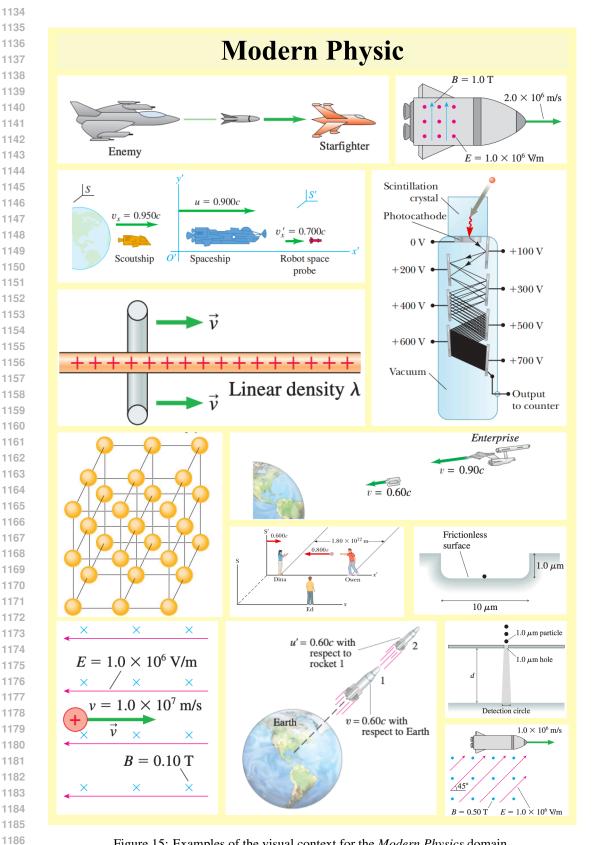


Figure 15: Examples of the visual context for the Modern Physics domain.

### **CoT Pro**r

### **CoT Prompting for Generating Answer**

Please answer the question with step by step reasoning.

Figure 16: CoT prompting for generating answer.

# Rule-based Answer Extraction (MC) def MetaPhyX\_process\_line\_MC(line): ret = {} answers = str(line['answer']) ret["index"] = line["index"] ret['gt'] = answers

ret['pred'] = line['prediction'].strip()

pattern = r'\b(?:correct|answer|option|Answer|Option|
Correct)\b[\s\S]\*?([A-D])'
 match = re.search(pattern, ret['pred'])

Figure 17: Rule-based answer extraction strategy for MC questions.

### G.1 COT PROMPTING FOR GENERATING ANSWER

The CoT prompting for generating answer is shown in Figure 16.

### G.2 RULE-BASED ANSWER EXTRACTION

The rule-based answer extraction strategies for MC and OE questions are shown in Figure 17 and Figure 18, respectively.

### G.3 PROMPT FOR ANSWER JUDGE

The prompt for answer judge is shown in Figure 19.

### G.4 PROMPT FOR CAPTION GENERATION

The prompt for caption generation is shown in Figure 20

### G.5 PROMPT FOR REASONING TYPE LABELING

The prompt for reasoning type labeling is shown in Figure 21 and Figure 22

1245 1246

1247

1248

1249 1250

1251

1252

1254

1256

1258

1260

12621263

1264

1265

1266

1267

1268

1269

1270

1271

1273

1277

1280 1281

1282

1283

1284

1285

1286

1288

1290

1292

1294

# Rule-based Answer Extraction (OE) def MetaPhyX\_process\_line(line): ret = {} answers = str(line['answer']) ret["index"] = line["index"] ret['gt'] = answers ret['pred'] = line['prediction'].strip() pattern = r'\b(?:final\s+answer| correct\s+answer)\b[^::]\*[::]\s\*(.\*?)(?=\n\n\n\\Z)' flags = re.IGNORECASE | re.DOTALL match = re.search(pattern, ret['pred'], flags=flags)

Figure 18: Rule-based answer extraction strategy for OE questions.

### **Prompt for Answer Judge**

```
Please read the following example.
Given predicted answer and ground truth answer, compare the these two
answers, then ONLY output judegement 1/0 for matched/unmatched at the
end of the prompt.
If the meaning is expressed in the same way, it is also considered
consistent, for example, 0.5m and 50cm.
If the given predicted mentions "approximately", then allow the
Approximation Error, such as 0.49 and approximately 0.5, 0.81 and
approximately 0.8.
Ground truth answer: 26.7kg \n
Predicted answer: The mass of block \( B \) is:
\boxed{26.7 \, \text\{kg\}}
\] \n
Judegement: 1
Ground truth answer: 46.3 kN \n
Predicted answer: The tension (T_B ) in the cable is approximately:
\boxed{46300 \, \text{N}}}
\1 \n
Judegement: 1
Ground truth answer: 12 m/s \n
Predicted answer: The speed of the box after 2.00 seconds is:
\boxed{11.3 \, \text{m/s}}
\] \n
Judegement: 0
Ground truth answer: 36.00 kg \n
Predicted answer: The mass of the hanging block ( m_2 ) must be
approximately:
\boxed{36.1 \, \text\{kg\}}
\] \n
Judegement: 1
Ground truth answer: 4.7 m \n
Predicted answer: The stuntman and villain slide approximately **4.69
meters**.
Judegement: 1
Ground truth answer: {}
Predicted answer: {}
Judegement:
```

Figure 19: Rule-based answer extraction strategy for OE questions.

# Prompt for Caption Generation Describe the fine-grained content of the image or figure, including scenes, objects, relationships, and any text present.

Figure 20: Prompt template for caption generation.

Physical Reasoning	Description
Physical Model Grounding Reasoning	This reasoning involves connecting the specific details of a problem description to fundamental physical concepts, laws, and idealized models. It's the process of identifying which area of physics is relevant and selecting the appropriate simplified representations that allow the problem to be analyzed using established physical principles and equations. Essentially, it translates a real-world or described scenario into a solvable physics framework.
Spatial Relation Reasoning	This focuses on understanding and manipulating the geometric and directional aspects of a physics problem. It involves visualizing the setup determining the positions, orientations, distances, angles, and relative movements of objects. This often requires using coordinate systems, vectors (including resolving them into components), and geometric principles.
Multi-Formula Reasoning	This reasoning type is required when a problem cannot be solved using a single physics equation. It involves identifying multiple relevant formulas or principles and understanding how they interrelate. The process typically involves using the output of one formula as the input for another, or setting up and solving a system of simultaneous equations derived from different physical laws.
Implicit Condition Reasoning	This involves recognizing and utilizing information or constraints that are not explicitly stated in the problem text but are implied by the context standard physics assumptions, or specific keywords. Examples include understanding that "starts from rest" means the initial velocity is zero, a "smooth" surface implies zero friction, a "light string" or "light pulley" means its mass is negligible, or that an object reaching its maximum height has a momentary vertical velocity of zero.
Numerical Reasoning	This reasoning refers to problems where solving requires the application of advanced mathematical methods beyond basic algebra and trigonometry. This includes techniques such as calculus, solving differential equations that model the system, vector calculus, Fourier analysis, linear algebra for complex systems, or other higher-level mathematical procedures necessary to manipulate the physical formulas and arrive at a solution. This applies when the mathematical technique itself is a core part of solving the physics regardless of whether the final answer is purely numerical or symbolic.
Predictive Reasoning	This involves using established physical laws and the initial conditions of a system to forecast its future state or behavior. Based on the principles governing the situation, you calculate or deduce what will happen after a certain time or interaction. Examples include predicting the trajectory of a projectile, the final temperature of a mixture after thermal equilibrium is reached, or the velocity of objects after a collision.

Table 6: Definitions of six physical reasoning categories in PHYX.

```
1350
                                                                                        Prompt for Reasoning Type Labeling(1)
1351
1352
                                                                                        You are an expert AI assistant specializing in analyzing physics
1353
1354
                                                                                        **Your lask:**
Your goal is to carefully read the provided physics problem and identify the **zero, one, or two MOST critical/dominant** reasoning types required to solve it, based on the definitions below. Your primary task is selection and prioritization. Assign a **maximum of two** labels per problem.
1355
1356
1357
                                                                                        **Reasoning Type Definitions:**
1358

    **Physical Model Grounding Reasoning:**
        *Explanation: Connecting problem details to physical concepts, laws, and idealized models (e.g., point mass, frictionless surface, ideal gas). Translating the scenario into a physics framework.

1359
1360
                                                                                       2. **Spatial Relation Reasoning:**
1361
                                                                                        \star Explanation: Understanding and manipulating geometric aspects (positions, angles, vectors, diagrams, coordinate systems).
1362
                                                                                       3. **Multi-Formula Reasoning:**
    * Explanation: Requiring the combination or sequential use of multiple distinct physics formulas or principles to find the solution.
1363
1364

    **Implicit Condition Reasoning:**
    * Explanation: Recognizing and using conditions not explicitly tated but implied by context or keywords (e.g., "starts from rest", smooth surface", "maximum height").

1365
1366
                                                                                        5. **Numerical Reasoning:** (Revised Definition)
                                                                                       1368
1369
                                                                                       6. **Predictive Reasoning:**
     *Explanation: Using physical laws and initial conditions to
forecast a future state or behavior (e.g., final velocity, trajectory,
final temperature).
1370
1371
1372
1373
1374
1375
                                                                                   Figure 21: Prompt for reasoning type labeling (1).
1376
1377
1378
                                                                                        Prompt for Reasoning Type Labeling(2)
1379
```

```
**Instructions:**
1. **Read and Analyze:** Carefully understand the problem and the likely steps/concepts needed for its solution.
2. **Identify Potential Types:** Determine which of the 6 reasoning types are involved in the solution process.
3. **Prioritize and Select:** From the potentially involved types, select **at most two** that are the **most critical, dominant, or uniquely challenging** aspects of solving this *particular* problem.
* Think about what makes the problem non-trivial. Is it complex geometry? Combining multiple physics laws? Recognizing hidden conditions? Needing calculus?

* If several types apply, choose the one or two that best represent the core difficulty or the essential nature of the solution process.
* If only one type truly stands out as the most essential characteristic, list only that one.

* If the problem is exceptionally simple or doesn't strongly fit any category as "most critical", output zero labels ('[']').
4. **Output Format:** **CRUTICAL:** Your entire response must consist *only* of a single Python-style list containing strings of the exact names for the selected zero, one, or two most critical reasoning types.

Oo **NOT** include any introductory text, explanations, labels, applogies, or any characters outside of the list itself.

***Correct Format Example (Two Types):** ["Reasoning Type A", "Reasoning Type B"]

***Correct Format Example (One Type):** ["Reasoning Type C"]

***Incorrect Format Example (One Type):** ["Reasoning Type C"]

***Example demonstrating the task (internal analysis, not part of the
                                   **Read and Analyze:** Carefully understand the problem and the
        **Example demonstrating the task (internal analysis, not part of the
    * **Problem:** "A 2 kg block, initially at rest on a frictionless horizontal surface, is pulled by a constant horizontal force of 10 N. what is its velocity after it has traveled 5 meters?"

* **(Internal Analysis:** Potential types involved: Grounding, Spatial (trivial here), Multi-formula (Fema then kinematics), Implicit ('at rest'), Predictive. Which are *most critical* (max 2)? Combining Fema and kinematics ('Multi-Formula Reasoning') is the core physics calculation. Recognizing 'at rest' ('Implicit Condition Reasoning') is crucial for the setup. These seem most central.)*

* **Required Output for this Example:** \""Multi-Formula Reasoning", "Implicit Condition Reasoning")
          **Now, analyze the following physics problem:**
            --- Problem Start ---
```

Figure 22: Prompt for reasoning type labeling (2).

### H CASE STUDY

### **List of Case Study Figures**

1	Mechanics 1: Correct Case	3
2	Mechanics 2: Correct Case	)
3	Mechanics 3: Visual Reasoning Error	)
4	Mechanics 4: Text Reasoning Error	l
5	Mechanics 5: Lack of Knowledge	)
6	Electromagnetism 1: Correct Case	;
7	Electromagnetism 2: Correct Case	ļ
8	Electromagnetism 3: Visual Reasoning Error	į
9	Electromagnetism 4: Text Reasoning Error	í
10	Electromagnetism 5: Lack of Knowledge	7
11	Thermodynamics 1: Correct Case	
12	Thermodynamics 2: Correct Case	)
13	Thermodynamics 3: Visual Reasoning Error	)
14	Thermodynamics 4: Text Reasoning Error	l
15	Thermodynamics 5: Lack of Knowledge	)
16	Wave/Acoustics 1: Correct Case	
	Wave/Acoustics 2: Correct Case	ļ
	Wave/Acoustics 3: Visual Reasoning Error	į
	Wave/Acoustics 4: Text Reasoning Error	í
20	Wave/Acoustics 5: Lack of Knowledge	-
21	Optics 1: Correct Case	
22	Optics 2: Correct Case	
23	Optics 3: Visual Reasoning Error	
24	Optics 4: Text Reasoning Error	
25	Optics 5: Lack of Knowledge	
26	Modern Physics 1: Correct Case	
27	Modern Physics 2: Correct Case	
28	Modern Physics 3: Visual Reasoning Error	
20	Modern Physics 4: Text Reasoning Error	
30	Modern Physics 5: Lack of Knowledge	
50	INTOGETH I HYDRED J. LACK OF KHOWIEUZE	/

Domain	Correct	Visual Reasoning Error	Text Reasoning Error	Lack of Knowledge
Mechanics	1, 2	3	4	5
Electromagnetism	6, 7	8	9	10
Thermodynamics	11, 12	13	14	15
Wave/Acoustics	16, 17	18	19	20
Optics	21, 22	23	24	25
Modern Physics	26, 27	28	29	30

Table 7: Table index of case study figures by domains with associated error categories.

### 1458 1459 **Mechanics** 1460 1461 1462 Description: 1463 A firefighter, a distance d from a burning building, directs a stream of water from a fire hose at angle $\theta_i$ above the 1464 horizontal as shown in figure. 1465 1466 Question: 1467 If the initial speed of the stream is vi, at what height h does 1468 the water strike the building? 1469 1470 1471 1472 GPT-40: 1473 To determine the height $\setminus$ ( h $\setminus$ ) at which the water strikes the building, we will analyze the motion of the water as a projectile. The motion can be broken into horizontal and vertical components. Here's the 1474 step-by-step reasoning: 1475 Step 1: Break the initial velocity into components The initial velocity $(v_i)$ can be resolved into horizontal and vertical components: $[v_i]$ = $v_i$ 1476 $\cos \theta_i \le v_i \le v_$ 1477 1478 Step 2: Time of flight to reach the building 1479 takes for the water to reach the building is given by: $\[ d = v_{ix} t \]$ 1480 Substitute $(v_{ix} = v_i \cos \theta_i): [t = \frac{d}{v_i \cos \theta_i}]$ 1481 Step 3: Vertical displacement at time $\setminus$ ( t $\setminus$ ) 1482 The vertical displacement $\setminus$ ( y $\setminus$ ) of the water at time $\setminus$ ( t $\setminus$ ) is determined by the kinematic equation: 1483 $[ y = v_{iy} t - \frac{1}{2} g t^2 ]$ Substitute $\langle v_{iy} = v_i \sin \theta_i \rangle$ and $\langle t = \frac{d}{v_i \cos \theta_i} \rangle$ : 1484 $\label{eq:cost} $$ y = \left( v_i \sin \theta_i \right) \left( \frac{d}{v_i \cos \theta_i} \right) - \frac{1}{2} g \left( \frac{d}{v_i} \cos \theta_i \right) - \frac{1}{2} g \left( \frac{d}{v_i} \right) - \frac{1}{2} g \left( \frac{$ 1485 $\cos \theta_i \$ 1486 Simplify the terms: $\label{eq:cos-2} $$ [y = \frac{d \tan \theta_i}{1} - \frac{g d^2}{2 v_i^2 \cos^2 \theta_i} ] $$$ 1487 1488 Step 4: Interpret the result The height $\sqrt{(h \lor)}$ at which the water strikes the building is the vertical displacement $\sqrt{(y \lor)}$ when the 1489 water reaches the horizontal distance \( d \). Thus: 1490 1491 Final Answer: 1492 The height $\setminus$ ( h $\setminus$ ) at which the water strikes the building is: 1493 $[ h = d \times_i - \frac{g d^2}{2 v_i^2 \cos^2 \theta_i} ]$ 1494 1495 - (d) is the horizontal distance to the building. 1496 - \(\theta\_i\) is the angle of the water stream above the horizontal, - $\ (v_i \ )$ is the initial speed of the water, 1497 - $\ (g \ )$ is the acceleration due to gravity ( $\ (9.8 \ , \text{text}{m/s}^2 \ )$ ). 1498

Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning

### Ground Truth: dtanθi-gd2/2vi2cos2θi

1499

1500 1501

15021503

1506

Figure 1: A sample correct case of Mechanics. Back to List of Figures | Back to Table Index

### **Mechanics**



### Description:

A landscape architect is designing an artificial waterfall in a city park. Water flows at a speed of 1.70 m/s as it exits the end of a horizontal channel located at the top of a vertical wall of height h=2.35 m. From there, the water falls freely into a pool below, as shown in the figure.

### Question:

To sell her plan to the city council, the architect wants to build a model to standard scale, which is one-twelfth actual size. How fast should the water flow in the channel in the model?

### Options:

A. 0.325m/s

B. 0.120m/s

C. 0.491m/s

D. 0.212m/s



### GPT-40:

The speed of water in the model can be determined using the principle of similarity. The velocity in the model is scaled by the square root of the length scale factor. Since the model is one-twelfth the size of the actual waterfall, the length scale factor is  $\$  (1/12  $\$ ).

The velocity in the model is given by:

 $\label{lem:condition} $$ v_{\text{\cod}} = v_{\text{\cod}} \cdot \left[ v_{\text{\cod}} \cdot \right] \$ 

Substituting the values:

 $\label{lem:condition} $$ v_{\text{model}} = 1.70 \cdot \sqrt{1}{12} = 1.70 \cdot \sqrt{0.0833} \cdot 1.70 \cdot 0.2887 \cdot 0.491 \cdot \cdot \sqrt{1.0833} \cdot 1.70 \cdot 0.0833 \cdot 0.0891 \cdot 0.089$ 

The correct answer is:  ${\color{red} {\bf C}}$ 

Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning

Ground Truth: C. 0.491m/s

Figure 2: A sample correct case of Mechanics. Back to List of Figures | Back to Table Index

### 1566 1567 Mechanics 1568 1569 1570 Description: 1571 We wrap a light, nonstretching cable around a solid 2.0 m 9.0 N cylinder, of mass 50 kg and diameter 0.120 m, that rotates in 1572 frictionless bearings about a stationary horizontal axis as 1573 shown in figure. We pull the free end of the cable with a 1574 constant 9.0 N force for a distance of 2.0 m; it turns the 50 kg 0.120 n cylinder as it unwinds without slipping. The cylinder is 1575 initially at rest. 1576 Question: 1577 Find the final speed of the cable. 1578 1579 1580 To solve this problem, we will use the principles of work and energy, as well as rotational dynamics. 1581 Here's the step-by-step reasoning: Step 1: Understand the system 1584 1585 Step 2: Work done by the force 1586 1587 Step 3: Rotational kinetic energy of the cylinder Step 4: Relationship between linear and angular velocity 1590 1591 Step 5: Total kinetic energy 1592 The total kinetic energy of the system is the sum of the rotational kinetic energy of the cylinder and the 1593 translational kinetic energy of the cable: 1594 \[ K\_{\text{total}} = K\_{\text{rot}} + K\_{\text{trans}}, \] where $\ (K_{\text{trans}}) = \frac{1}{2} \text{ m v}^2 \ )$ 1595 Substitute $\ (K_{\text{rot}} = \frac{1}{2} I \geq^2 \)$ and $\ (\omega = \frac{v}{r} \)$ : 1596 \[ K\_{\text{rot}}} {0.060}\right)^2. \] 1597 Simplify: 1598 $[K_{\text{tot}}] = \frac{1}{2} \cdot 0.090 \cdot \frac{rac{v^2}[0.060^2]} = \frac{0.090}{2} \cdot 0.060^2 \cdot v^2 = 12.5$ v^2. \1 The translational kinetic energy is: $[K_{\text{trans}}] = \frac{1}{2} \text{ m v}^2 = \frac{1}{2} (50) \text{ v}^2 = 25 \text{ v}^2.$ 1601 Thus, the total kinetic energy is: $[K_{\text{text{total}}} = K_{\text{text{rot}}} + K_{\text{text{trans}}} = 12.5 \text{ v}^2 + 25 \text{ v}^2 = 37.5 \text{ v}^2.$ 1602 Step 6: Work-energy principle 1604 Final Answer: 1606 The final speed of the cable is approximately: \[ \boxed{0.69 \, \text{m/s}}. \] 1608 Reasoning Type: Multi-Formula Reasoning, Physical Model Grounding Reasoning 1609 1610 **Error Type: Text Reasoning Errors** 1611 1612 Ground Truth: 1.2 m/s 1613 1614

Figure 3: A sample error case of Mechanics. Error category: Visual Reasoning Error Back to List of Figures | Back to Table Index

1615

### 1620 1621 **Mechanics** 1622 1623 1624 Description: 1625 We want to slide a 12-kg crate up a 2.5-m-long ramp inclined at 30°. A worker, ignoring friction, calculates that Point 1626 2 he can do this by giving it an initial speed of 5.0 m/s at the 1627 bottom and letting it go. But friction is not negligible; the 1628 crate slides only 1.6 m up the ramp, stops, and slides back down. 1629 1630 Question: 1631 How fast is the crate moving when it reaches the bottom of Point (1), (3) the ramp? 1632 1633 1634 GPT-40: To solve this problem, we will use the principles of energy conservation, taking into account the work 1635 done by friction. Here's the step-by-step reasoning: 1636 1637 Step 1: Define the system and energy components 1638 1639 Step 2: Write the energy conservation equation 1640 1641 Step 3: Calculate the work done by friction\*\* 1642 The work done by friction is: $\[ W_{\text{text\{friction\}}} = f_k d \]$ 1643 The frictional force $\ (f_k \ )$ is given by: $\ [f_k = \mu_k \ ]$ 1644 \(\mu\_k\): Coefficient of kinetic friction. 1645 The total distance traveled by the crate is: 1646 $[d = 1.6 \, \text{(up)} + 1.6 \, \text{(down)} = 3.2 \, \text{(m)}. ]$ 1647 Thus: $\[ W_{\text{friction}} = \mu_k m g \cos \theta \]$ 1648 Step 4: Solve for the final velocity 1649 1650 Step 5: Substitute known values 1651 Given: 1652 $(v_1 = 5.0 \setminus, \text{text}\{m/s\} \setminus),$ - $\ (g = 9.8 \ , \text{text}{m/s}^2 \ )$ 1653 - \( \theta = 30^\circ \), - $\ (d = 3.2 \ , \text{text{m}} \ )$ 1655 - \( \mu\_k = \text{unknown (not provided)} \). 1656 1657 Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning 1658 1659 **Error Type: Text Reasoning Errors** 1660 1661 Ground Truth: 2.5 m/s 1662

Figure 4: A sample error case of Mechanics. Error category: Text Reasoning Error Back to List of Figures | Back to Table Index

1663 1664

1665

### 1674 1675 **Mechanics** 1676 1677 1678 Description: 25.0 1679 A 500.0 g bird is flying horizontally at 2.25 m/s, not paying Bird cm much attention, when it suddenly flies into a stationary 1680 vertical bar, hitting it 25.0 cm below the top as shown in 1681 figure. The bar is uniform, 0.750 m long, has a mass of 1.50 1682 kg, and is hinged at its base. 1683 1684 Question: 1685 What is the angular velocity of the bar just as it reaches the ground? Hinge 1687 1688 GPT-40: To solve this problem, we need to determine the angular velocity of the bar just as it reaches the ground. 1689 This involves analyzing the collision and the subsequent rotational motion of the bar. Here's the step-bystep solution: Step 1: Conservation of Angular Momentum 1692 The bird collides with the bar, imparting angular momentum to it. Since the hinge is frictionless, 1693 angular momentum about the hinge is conserved during the collision. 1694 Angular Momentum Before Collision The bird has linear momentum before the collision: 1695 $[p_{\text{text}}] = m_{\text{text}} v_{\text{bird}} \]$ The bird's angular momentum about the hinge is: 1698 $\ L_{\text{bird}} = p_{\text{bird}} \ cdot \ r \ ]$ 1699 where $\langle (r = 0.750 - 0.250 = 0.500 \setminus, \text{text}\{m\} \rangle)$ is the perpendicular distance from the hinge to the point of collision. 1700 $[L_{\text{text}}] = (1.125)(0.500) = 0.5625 \, \text{kg·m}^2/\text{text}[s] \]$ 1701 Angular Momentum After Collision After the collision, the bar rotates about the hinge. Its angular momentum is: 1702 $\ [L_{\text{bar}} = I_{\text{bar}} \setminus \$ 1703 1704 - \( I\_{\text{bar}} \) is the moment of inertia of the bar about the hinge, - \(\omega\) is the angular velocity of the bar. 1705 1706 Using conservation of angular momentum: \[ L\_{\text{bird}} = L\_{\text{bar}} \] \[ 0.5625 = (0.28125) \cdot \omega \] 1708 Solve for \( \omega \): 1709 $\[ \omega = \frac{0.5625}{0.28125} = 2.00 \, \text{$\ \]} \]$ Thus, the angular velocity of the bar immediately after the collision is: 1710 \[ \omega = 2.00 \, \text{rad/s} \] 1711 1712 Step 2: Energy Conservation During Rotation After the collision, the bar rotates about the hinge. As it falls, its potential energy is converted into 1713 rotational kinetic energy. 1714 Initial Potential Energy 1715 Solve for \(\omega\): 1716 \[\omega^2 = \frac{2(5.5125)}{0.28125} = 39.2 \] 1717 \[ \omega = \sqrt{39.2} = 6.26 \, \text{rad/s} \] 1718 Final Answer: 1719 The angular velocity of the bar just as it reaches the ground is: \[ \boxed{6.26 \, \text{rad/s}} \] 1720 1721 Reasoning Type: Multi-Formula Reasoning, Physical Model Grounding Reasoning 1722 1723 Error Type: Lack of knowledge 1724 1725 Ground Truth: 6.58 rad/s 1726

Figure 5: A sample error case of Mechanics. Error category: Lack of Knowledge Back to List of Figures | Back to Table Index

### 1728 1729 1730 1731 1732 1733 1734 **Electromagnetism** 1735 1736 1737 Description: 1738 In the slidewire, energy is dissipated in the circuit owing to 1739 its resistance. Let the resistance of the circuit (made up of the moving slidewire and the U-shaped conductor that 1740 connects the ends of the slidewire) at a given point in the 1741 slidewire's motion be R. 1742 Question: 1743 Find the rate at which energy is dissipated in the circuit. 1744 1745 1746 1747 GPT-40: 1748 To find the rate at which energy is dissipated in the circuit, we need to calculate the power dissipated 1749 due to the resistance $\setminus$ ( R $\setminus$ ). The power dissipated in a resistor is given by: $[P = \frac{Fac}{mathcal{E}^2}{R}, ]$ 1750 where \(\) \mathcal{E} \\) is the electromotive force (emf) induced in the circuit. Let's calculate 1751 \(\mathcal{E}\) and then substitute it into the formula for power dissipation. 1752 Step 1: Determine the induced emf (\( \mathcal{E} \)) 1753 1754 According to Faraday's law of electromagnetic induction, the emf is given by: $\[ \]$ $\[ \]$ $\[ \]$ B L v, $\]$ where: 1755 $\sim$ ( B $\sim$ ) is the magnitude of the uniform magnetic field (perpendicular to the plane of the circuit), 1756 - \( L \) is the length of the slidewire (the part of the circuit moving in the magnetic field), - $\setminus$ (v $\setminus$ ) is the velocity of the slidewire (perpendicular to both $\setminus$ ( $\setminus$ mathbf{B} $\setminus$ ) and $\setminus$ ( L $\setminus$ )). 1757 1758 Step 2: Calculate the current in the circuit 1759 The current $\setminus$ ( I $\setminus$ ) in the circuit is related to the emf and the resistance $\setminus$ ( R $\setminus$ ) by Ohm's law: $[I = \frac{\mathbb{E}}{R}. ]$ 1760 1761 Step 3: Calculate the power dissipated in the circuit 1762 The power dissipated in the circuit due to the resistance $\ (R )$ is given by: $\ P = I^2 R. \]$ 1763 Substitute $\ (I = \frac{B L v}{R} \)$ into this expression: $\ P = \left(\frac{B L v}{R} \right)^2 R. \]$ 1764 Simplify the expression: $\ [P = \frac{(B L v)^2}{R}. \]$ 1765 Final Answer: 1766 The rate at which energy is dissipated in the circuit is: 1767 $[P = \frac{(B L v)^2}{R}. ]$ 1768 This is the power dissipated in the circuit due to its resistance. 1769 1770 Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning 1771

Figure 6: A sample correct case of Electromagnetism.

Back to List of Figures | Back to Table Index

1772

1773 1774 1775

1776

1777 1778

17801781

Ground Truth: B<sup>2</sup>L<sup>2</sup>v<sup>2</sup>/R

### 1782 1783 1784 1785 Electromagnetism 1786 1787 1788 Description: Question: 1789 A sensitive electronic device with a resistance What is the R-L time constant $\tau$ ? 1790 of R=175 $\Omega$ is to be connected to a source of emf (with negligible internal resistance) by a switch. 1791 The device is designed to operate with a 1792 current of 36 mA, but to avoid damage to the device, the current can rise to no more than 1793 4.9mA in the first 58 µs after the switch is 1794 closed. An inductor is therefore connected in 1795 series with the device, as shown in the figure; the switch in question is S1. 1796 1797 Options: 1798 Α. 220μΗ B. 390µH C. 750µH D. 280µH Closing switch S1 connects the R-L combination 1800 1801 1802 1803 1805 1806 Closing switch $S_2$ while opening switch $S_1$ 1807 1808 GPT-40: 1809 1810 $\[ \text{tau} = \frac{L}{R}. \]$ 1811 The current in an \( R\text{-}L \) circuit grows according to: 1812 $[i(t) = I_{\text{max}} \left[ 1 - e^{-t/\lambda u} \right], ]$ where $\ (I_{\max}) = 36 \ \mathrm{mathrm\{mA\}} \ )$ , $\ (i(t) = 4.9 \ \mathrm{mathrm\{mA\}} \ )$ , and $\ (t = 58 \ \mathrm{mathrm\{mA\}} \ )$ 1813 $\mu\mathrm{s}\).$ 1815 $[e^{-t/\lambda u} = 1 - \frac{i(t)}{I_{\lambda u}} = 1$ 1816 1817 Substitute the values: \[ e^{-58 \times 10^{-6}}\tau\} = 1 - \frac{4.9 \times 10^{-3}}{36 \times 10^{-3}} = 1 - 0.1361 = 0.8639. \] 1818 1819 Take the natural logarithm: \[ -\frac{58 \times 10^{-6}}{\tau} = \ln(0.8639). \] 1820 $[ \tau = -\frac{58 \times 10^{-6}}{\ln(0.8639)}. ]$ 1821 \[ \ln(0.8639) \approx -0.146. \] 1822 \[\tau = \frac{58 \times 10^{-6}}{0.146} \approx 397 \times 10^{-6}\ \mathrm{s}. \] Now, calculate $\ \ L \ )$ : 1824 \[ L =\tau R = (397 \times 10^{-6})(175) \approx 69.5 \times 10^{-3}\ \mathrm{H} = 390 \ \mu\mathrm{H}.\] 1825 Correct option: B 1826 1827 Reasoning Type: Multi-Formula Reasoning, Numerical Reasoning 1828 1829 Ground Truth: B. 390µH 1830 1831

Figure 7: A sample correct case of Electromagnetism.

Back to List of Figures | Back to Table Index

1833

### Electromagnetism Description: In the circuit, the switch S has been open for a long time and is suddenly closed. Neither the battery nor the 4Rinductors have any appreciable resistance. 15 V Question: What does the voltmeter read 0.115 ms after S is closed? To solve this problem, we need to analyze the circuit and determine the voltage across the voltmeter \( V \) 0.115 ms after the switch \( S \) is closed. Here's the step-by-step reasoning: ### Step 1: Understand the circuit - The circuit consists of a 20.0 V battery, resistors, and inductors. The switch \( S \) has been open for a long time, so initially, no current flows through the circuit. - When \( S \) is closed, current starts flowing, and the inductors oppose the change in current due to their inductance. Step 2: Inductors' behavior Step 3: Equivalent inductance of the parallel inductors Step 4: Time constant of the RL circuit Step 5: Current in the circuit as a function of time Step 6: Voltage across the parallel inductors $\ \ V = L_{\text{eq}} \ frac{dI}{dt} \ ]$ $Substitute \ (L_{\text{eq}} = 4.87 \ , \text{text}\{mH\} = 4.87 \ times \ 10^{-3} \ , \text{text}\{H\} \ ) \ and \ ( \ frac\{dI\}\{dt\} = 699 \ ), \ ( \ frac\{dI\}\{dt\} = 699 \$ \text{A/s} \): $[V = 4.87 \times 10^{-3} \cdot 699 \cdot 3.41 , \text{$\ \]}$ Final Answer: The voltmeter reads approximately \*\*3.41 V\*\* 0.115 ms after the switch is closed. Reasoning Type: Physical Model Grounding Reasoning, Numerical Reasoning **Error Type: Visual Reasoning Error** Ground Truth: 9.0 V

Figure 8: A sample error case of Electromagnetism. Error category: Visual Reasoning Error Back to List of Figures | Back to Table Index

### Electromagnetism Description: Figure shows an electric dipole in a uniform electric field of magnitude 5.0 \* 10^5 N/C that is directed parallel to the É plane of the figure. The charges are pm 1.6 \* 10^-19 C; both lie in the plane and are separated by 0.125 \* 10^-9m. Question: Find the magnitude of the torque. To calculate the torque (\(\\\\\\\\)) acting on the electric dipole in the uniform electric field, we use the formula:\[\tau = pE \sin \theta\] where: - $\setminus$ ( p $\setminus$ ) is the dipole moment, - \( E \) is the magnitude of the electric field, $- (\theta \ )$ is the angle between the dipole moment $(\theta \ )$ and the electric field $(\theta \ )$ . Step 1: Calculate the dipole moment (\( ( p \)) The dipole moment is given by: [p = qd]where: Substitute the values: $\label{eq:p = (1.6 \times 10^{-19})(0.125)} $$ [p = 2.0 \times 10^{-20} \times Cm] $$ ]$ Step 2: Use the torque formula Final Answer: The magnitude of the torque is: Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning **Error Type: Text Reasoning Errors** Ground Truth: 5.7 \* 10^{-24} \, \text{N} \cdot \text{m}

Figure 9: A sample error case of Electromagnetism. Error category: Text Reasoning Error Back to List of Figures | Back to Table Index

#### **Electromagnetism** Description: Point charges $q_1 = +12$ nC and $q_2 = -12$ nC are 0.100 m apart. (Such pairs of point charges with equal magnitude and opposite sign are called electric dipoles. Question: Compute the total field at point c. To compute the total electric field at point (c) due to the charges $(q_1)$ and $(q_2)$ , we will calculate the contributions of the electric fields from each charge and then combine them vectorially. Here's the step-by-step reasoning: Step 1: Understand the setup $- (q_1 = +12 \setminus \text{text}(nC) \setminus \text{positive charge}).$ $- (q_2 = -12 \setminus \text{text}\{nC\} \setminus \text{negative charge}).$ - Distance between $\ \ (q_1 \ ) \ and \ \ (q_2 \ ): \ \ (0.100 \ , \ \text{text}\{m\} \ ).$ - The electric field at $\langle (c \rangle)$ due to each charge is directed along the line joining the charge and point $\langle (c \rangle)$ **\).** Step 2: Electric field due to a point charge The electric field due to a point charge is given by: $[E = \frac{k|q|}{r^2}, ]$ where: $- (k = 8.99 \times 10^9 , \text{N·m}^2/\text{text}(C)^2 ) (Coulomb's constant),$ - \( q \) is the magnitude of the charge, $- \setminus (r \setminus)$ is the distance from the charge to the point where the field is being calculated. Step 3: Calculate the magnitude of the electric field at $\ (c \ )$ due to $\ (q_1 \ )$ and $\ (q_2 \ )$ \*\* Final Answer: The total electric field at point $\setminus$ ( c $\setminus$ ) is: $\[ \boxed{11.8 \, \text{N/C} \, \text{upward along the vertical axis)}}. \]$ Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning Error Type: Lack of Knowledge Ground Truth: 5.7 \* 10^{-24} N/C

Figure 10: A sample error case of Electromagnetism. Error category: Lack of Knowledge Back to List of Figures | Back to Table Index

#### 1998 1999 2000 Thermodynamics 2001 2002 2003 Description: 2004 A partition divides a thermally insulated box into two 2Vcompartments, each of volume V(figure). Initially, one 2005 compartment contains n moles of an ideal gas at 2006 temperature T, and the other compartment is evacuated. 2007 We break the partition and the gas expands, filling both Vcompartments. 2008 2009 **Ouestion:** What is the entropy change in this free-expansion process? 2010 2011 2012 2013 GPT-40: 2014 To calculate the root-mean-square (rms) speed \(v\_{\mathrm{rms}}\)\) of the molecules in the gas, we use the formula: $\ [v_{\mathrm{rms}} = \sqrt{\frac{1}{N} \sum_{i=1}^N v_i^2} ]$ 2015 where: 2016 \(N\) is the number of molecules, 2017 - $\langle v_i \rangle$ is the speed of the $\langle i \rangle$ -th molecule, and $- \v_i^2 = v_{ix}^2 + v_{iy}^2 \), where \v_{ix}\) and \v_{iy}\) are the \xv_i and \y\)-components of the \xi_{ix}\) are the \xi_{ix}\)- and \xi_{ix}\)- and$ 2018 velocity of the \(i\)-th molecule. 2019 Step 1: Extract the velocity components for each molecule 2020 From the diagram, the velocity components for each molecule are: 2021 1. Molecule 1: $( \text{vec}\{v\}_1 = 10 \text{hat}\{i\} - 10 \text{hat}\{j\} \) \ (v_{1x} = 10, \, v_{1y} = -10 \)$ 2022 2. Molecule 2: $( \text{vec}\{v\}_2 = 2 \text{hat}\{i\} + 15 \text{hat}\{j\} ) (v_{2x} = 2, \cdot, v_{2y} = 15)$ 2023 6. Molecule 6: $( \text{vec}\{v\}_6 = 0 \text{hat}\{i\} - 14 \text{hat}\{j\} ) (v_{6x} = 0, \cdot, v_{6y} = -14)$ 2024 Step 2: Calculate the speed \(v\_i\) for each molecule 2025 The speed of each molecule is given by: $[v_i = \sqrt{v_{ix}^2 + v_{iy}^2}]$ 2026 1. Molecule 1: $[v_1 = \sqrt{10^2 + (-10)^2} = \sqrt{100 + 100} = \sqrt{200} = 10\sqrt{2} ]$ 2027 2. Molecule 2: $[v_2 = \sqrt{2^2 + 15^2} = \sqrt{4 + 225} = \sqrt{229}]$ 2028 6. Molecule 6: $[v_6 = \sqrt{0^2 + (-14)^2} = \sqrt{0 + 196} = 14]$ 2029 Step 3: Calculate \(v\_i^2\) for each molecule 2030 We now square the speeds: 2031 1. $(v_1^2 = (10 \setminus sqrt\{2\})^2 = 200)$ 2032 2. $(v_2^2 = (\sqrt{229})^2 = 229)$ 2033 6. $(v_6^2 = 14^2 = 196)$ 2034 2035 Step 4: Calculate the mean of \(v\_i^2\) The mean of $(v_i^2)$ is: 2036 $[\text{text}] = \frac{1}{N} \sum_{i=1}^N v_i^2 = \frac{1}{6} (200 + 229 + 100 + 104 + 61 + 196)$ First, calculate the sum: \[ 200 + 229 + 100 + 104 + 61 + 196 = 890 \] 2037 Now divide by (N = 6): $[ \text{Mean} = \frac{890}{6} \cdot 148.33 ]$ 2038 2039 Step 5: Calculate \(v\_{\mathrm{rms}}\) 2040 Finally, take the square root of the mean: $[v_{\mathrm{s}}] = \sqrt{148.33} \exp 12.18$ 2041 2042 Final Answer: The rms speed is approximately: $[v_{\text{mathrm{rms}}} \ge 12.2 \, \text{w.}]$ 2043 2044 Reasoning Type: Physical Model Grounding Reasoning, Numerical Reasoning 2045 2046 Ground Truth: 12.2m/s 2047

Figure 11: A sample correct case of Thermodynamics.

Back to List of Figures | Back to Table Index

2048 2049

# **Thermodynamics**



#### Description:

# Two piston/cylinder arrangements, A and B, have their gas chambers connected by a pipe, as shown in the figure. The cross-sectional areas of the pistons are $A_A=75~cm^2$ and $A_B=25~cm^2$ , with the piston mass in arrangement A being $m_A=25~kg$ . Assume an outside pressure of 100kPa and standard gravitational acceleration.

#### Question:

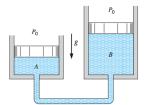
Find the mass  $m_B$  such that none of the pistons rests on the bottom.

#### Options:

A. 490kg B. 8.33kg

C. 154kg

D. 10.2kg



#### GPT-40

To solve this, we use the principle of pressure equilibrium. The pressure in both chambers must be equal since they are connected by a pipe.

Step 1: Pressure in chamber A

The pressure in chamber A is given by:

 $[P_A = P_0 + frac\{m_A g\}\{A_A\}]$ 

Substitute the values:

 $[P_A = 100,000 + \frac{25 \cdot 00075} = 100,000 + 32,700 = 132,700 , \text{}]$ 

Step 2: Pressure in chamber B

The pressure in chamber B is given by:

 $\[ P_B = P_0 + \frac{m_B g}{A_B} \]$ 

Since \(P\_A = P\_B\), we equate the two pressures: \[ 132,700 = 100,000 + \frac{m\_B \cdot 9.81}{0.0025} \]

Step 3: Solve for \(m\_B\)

Rearrange to isolate  $\mbox{(m_B\)}$ :

 $[\frac{m_B \cdot 0.0025} = 32,700]$ 

 $[ m_B = \frac{32,700 \cdot 0.0025}{9.81} = 8.33 \cdot, \text{kg} ]$ 

Final Answer: B

Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning

Ground Truth: B. 8.33kg

Figure 12: A sample correct case of Thermodynamics.

Back to List of Figures | Back to Table Index

#### 2106 2107 2108 2109 2110 **Thermodynamics** 2111 2112 2113 Description: 2114 The graph in figure shows a PV-diagram of the air in a human p (mm of Hg) 2115 lung when a person is inhaling and then exhaling a deep breath. Such graphs, obtained in clinical practice, are normally 12.0 2116 Maling somewhat curved, but we have modeled one as a set of 10.0 2117 straight lines of the same general shape. (Important: The 8.0 2118 pressure shown is the gauge pressure, not the absolute 6.0 pressure.) The process illustrated here is somewhat different 2119 from those we have been studying, because the pressure 4.0 2120 change is due to changes in the amount of gas in the lung, not 2.0 to temperature changes. Exha 2121 V(L) $\overline{o}$ 0.5 1.5 2122 Question: 2123 If the temperature of the air in the lung remains a reasonable 20°C, what is the maximum number of moles in this person's 2124 lung during a breath? 2125 2126 GPT-40: 2127 To determine the maximum number of moles of air in the lung during a breath, we can use the ideal gas 2128 $law: \setminus [pV = nRT \setminus]$ 2129 Where: 2130 $- \setminus (p \setminus)$ is the pressure (in pascals, Pa), 2131 $- \setminus (V \setminus)$ is the volume (in cubic meters, m<sup>3</sup>), - (n) is the number of moles, 2132 - (R) is the ideal gas constant (\( R = 8.314 \, \text{J/mol·K} \)), 2133 - $\setminus$ ( T $\setminus$ ) is the temperature (in kelvins, K). 2134 2135 Step 1: Identify the maximum pressure and volume 2136 From the graph: The maximum \*\*gauge pressure\*\* is $\ (p_{\text{auge}}) = 12.0 \ , \text{text}\{mmHg\} \ )$ . 2137 The maximum volume is $\ \ V = 1.5 \ , \ \text{text}\{L\} \ )$ . 2138 2139 Step 2: Convert pressure and volume to SI units 2140 2141 Step 3: Convert temperature to kelvins 2142 2143 Step 4: Solve for the number of moles 2144 2145 Final Answer: 2146 The maximum number of moles of air in the lung during a breath is approximately: \[\boxed{0.0634 \, \text{mol}}\] 2147 2148 Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning 2149 2150 **Error Type: Visual Reasoning Errors** 2151 2152 Ground Truth: 0.059mol 2153 2154

Figure 13: A sample error case of Thermodynamics. Error category: Visual Reasoning Error Back to List of Figures | Back to Table Index

2155

2156

#### Thermodynamics Description: A large tank of water has a hose connected to it. The tank is sealed at the top and has compressed air between the water surface and the top. When the water height h has the value 3.50 m, the absolute pressure p of the 4.00 m compressed air is $4.20 \times 10^5$ Pa. Assume that the air above the water expands at constant temperature, and take the atmospheric pressure to be 1.00 × 10^5 Pa. Question: As water flows out of the tank, h decreases. At what value 1.00 m of h does the flow stop? GPT-40: To determine the value of $\setminus$ (h $\setminus$ ) at which the flow stops, we need to analyze the system step by step. The flow stops when the pressure at the outlet of the tank equals the atmospheric pressure. Here's the Step 1: Understand the system Step 2: Pressure at the outlet Step 3: Relating the air pressure $\setminus$ ( p $\setminus$ ) to the water height $\setminus$ ( h $\setminus$ ) Step 4: Solve for \( h \) when the flow stops Substitute $\langle (p = \frac{4.20 \times 10^5) (0.50)}{4.00 - h} \rangle$ into the condition $\langle (p + \frac{4.20 \times 10^5) (0.50)}{4.00 - h} \rangle$ p\_{\text{atm}} \): \[\frac{(4.20 \times 10^5) (0.50)}{4.00 - h} + (1000)(9.8) h = 1.00 \times 10^5.\] \[\frac{2.10 \times 10^5}{4.00 - h} + 9800 h = 1.00 \times 10^5.\] \[\frac{2.10 \times 10^5}{4.00 - h} = 1.00 \times 10^5 - 9800 h.\] Step 5: Solve the quadratic equation Final Answer: The flow stops when $\ \ \ \ \ = 1.54 \ \ \ \ \ \ \ \ \ \ \ )$ . Reasoning Type: Multi-Formula Reasoning, Physical Model Grounding Reasoning **Error Type: Text Reasoning Errors** Ground Truth: 1.74m

Figure 14: A sample error case of Thermodynamics. Error category: Text Reasoning Error Back to List of Figures | Back to Table Index

# **Thermodynamics**

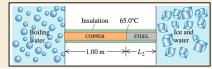


#### Description:

A long rod, insulated to prevent heat loss along its sides, is in perfect thermal contact with boiling water (at atmospheric pressure) at one end and with an ice-water mixture at the other. The rod consists of a 1.00 m section of copper (one end in boiling water) joined end to end to a length  $L_2$  of steel (one end in the ice-water mixture). Both sections of the rod have cross-sectional areas of 4.0 cm<sup>2</sup>. The temperature of the copper-steel junction is  $65.0^{\c}$  after a steady state has been set up.

#### Question:

What is the length  $L_2$  of the steel section?



#### GPT-40:

To solve for the length  $\ (L_2 \)$  of the steel section, we use the concept of steady-state heat conduction. In steady-state conditions, the heat flow rate through the copper section is equal to the heat flow rate through the steel section. This is because no heat is lost along the sides of the rod, and the system is in thermal equilibrium.

```
\label{thm:linear} Step 1: Write the heat flow rate for copper and steel For the copper section: $$ \left[ \frac{Cu}{t} = \frac{k_{\text{cu}} A (T_{\text{boiling}} - T_{\text{text{junction}}})}{L_{\text{text{Cu}}}, } \right] $$ ... $$ Step 2: Set the heat flow rates equal ... $$
```

#### Step 4: Substitute known values

- \( k\_{\text{Cu}} = 390\,\text{W/m·K} \), - \( k\_{\text{steel}} = 46\,\text{W/m·K} \),
- \( T\_{\text{boiling}} = 100.0^\circ\text{C} \),
- \( T\_{\text{junction}} = 65.0^\circ\text{C} \),
  - \( T\_{\text{ice}} = 0.0^\circ\text{C} \),
- \( L\_{\text{Cu}} = 1.00\,\text{m} \).

Substitute these values into the equation:

Step 5: Perform the calculations

Final Answer:

The length of the steel section is: \[ \boxed{0.219\,\text{m}}. \]

Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning

Error Type: Lack of Knowledge

Ground Truth: 5.7 \* 10^{-24} N/C

Figure 15: A sample error case of Thermodynamics. Error category: Lack of Knowledge Back to List of Figures | Back to Table Index

2268 2270 2271 2272 Wave/Acoustics 2273 2274 Description: 2276 A directional loudspeaker directs a sound wave of 5λ/4 2277 wavelength λ at a wall (figure) 2278  $3\lambda/4$ 2279 λ/4 Question: 2280 At what distances from the wall could you stand AN 2281 and hear no sound at all? 2282 Speaker 2285 GPT-40: To determine the distances from the wall where no sound is heard, we need to analyze the interference 2287 of the sound waves. Here's the step-by-step reasoning: Step 1: Understand the setup 2289 - A speaker emits sound waves toward a wall. - The sound waves reflect off the wall and interfere with the incoming waves. 2290 - At certain points, destructive interference occurs, resulting in no sound being heard. 2291 Step 2: Condition for destructive interference 2292 Destructive interference occurs when the path difference between the incident wave and the reflected 2293 wave is an odd multiple of half the wavelength (\(\ \\ \\ \): \[\ \text{Path difference}\] = (2m + 1)2294  $\frac{1}{2}$ ,  $quad m = 0, 1, 2, \\dots ]$ Here,  $\setminus$  ( m  $\setminus$ ) is an integer representing the order of the interference. 2295 2296 Step 3: Path difference in terms of distance from the wall - Let the distance from the wall be \( d \). 2297 The sound wave travels a distance  $\setminus$  ( d  $\setminus$ ) to the wall and another distance  $\setminus$  ( d  $\setminus$ ) back after reflection. 2298 - The total path difference is \( 2d \). 2299 For destructive interference:  $\ [2d = (2m + 1) \frac{2}{1}]$ 2300 Step 4: Solve for \( d \) 2301 2302 This gives the distances from the wall where destructive interference occurs. Step 5: List the distances For different values of  $\ (m \ )$ : - \( m = 0 \): \( d = \frac{\lambda}{4} \) 2305 -  $\ (m = 1 \): \ (d = \frac{3\lambda}{4} \)$ 2306  $- (m = 2 ): (d = \frac{5\lambda}{4} )$ 2307 - And so on... 2308 2309 Final Answer: The distances from the wall where no sound is heard are: 2310  $\label{lem:lembda} $$ d = \frac{14}, \frac{3\lambda}{4}, \frac{5\lambda}{4}, \frac{7\lambda}{4}, \frac{3\lambda}{4}, \frac$ 2311 2312 Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning

Figure 16: A sample correct case of Wave/Acoustics.

Back to List of Figures | Back to Table Index

23132314

2315 2316

2317

231823192320

Ground Truth:  $\lambda/4$ ,  $3\lambda/4$ ,  $5\lambda/4$ 

#### 2322 2323 2324 2325 2326 2327 2328 2329 2330 Wave/Acoustics 2331 2332 2333 Description: Question: 2334 Small speakers A and B are driven in phase at At what distance d will the sound from the speakers 2335 725Hz by the same audio oscillator. Both first produce destructive interference at the 2336 speakers start out 4.50m from the listener, but listener's location? 2337 speaker A is slowly moved away. 2338 Options: 2339 A. 0.237m B. 0.336m C. 0.569m 2340 2341 4.50 m 2342 2343 2344 2345 GPT-40: 2346 To determine the distance \( (d \) at which destructive interference first occurs, we use the condition for destructive interference: 2347 $\[ \Delta L = \frac{\lambda}{2}, \]$ 2348 where $\setminus (\Delta L \setminus )$ is the path difference between the two speakers, and $\setminus (\Delta L \setminus )$ is the wavelength 2349 of the sound. 2350 Step 1: Calculate the wavelength \(\lambda \) 2351 The speed of sound in air is approximately $(v = 343 , \text{text} \{m/s\})$ , and the frequency of the sound is $(f = 343 , \text{text} \{m/s\})$ 2352 = 725 \, \text{Hz} \). The wavelength is given by: $[\lambda = \frac{v}{f} = \frac{343}{725} \quad 0.473 \, \text{} ]$ 2353 2354 Step 2: Path difference for destructive interference For the first instance of destructive interference, the path difference $\setminus$ ( $\setminus$ Delta L $\setminus$ ) is: 2355 $[\Delta L = \frac{1}{2} = \frac{0.473}{2} \cdot 0.237 , \text{ } ]$ 2356 2357 The path difference $\setminus$ \ Delta L $\setminus$ \) is caused by moving speaker $\setminus$ \( A $\setminus$ \) a distance $\setminus$ \( d $\setminus$ \) farther from the 2358 listener. Thus: $\[ \Delta L = d. \]$ 2359 So: $\setminus$ [ d = 0.237 $\setminus$ , $\setminus$ text{m}. $\setminus$ ] 2360 Final Answer: A 2361 2362 Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning 2363 2364 Ground Truth: A. 0.237m

2365 2366

2367

Figure 17: A sample correct case of Wave/Acoustics. Back to List of Figures | Back to Table Index

D. 0.414m

2422

2423 2424

2426242724282429

2378 2379 2380 2381 2382 2383 Wave/Acoustics 2384 2385 2386 Description: 2387 Very finely ground particles of cork are sprinkled in the Piston Glass Piles of cork 2388 bottom of the tube before the piston is inserted. As the tube particles 2389 vibrating piston is slowly moved forward, there are a few positions that cause the cork particles to collect in 2390 small, regularly spaced piles along the bottom. Figure 400 Hz ◀ 2391 shows an experiment in which the tube is filled with 2392 pure oxygen and the piston is driven at 400 Hz. 123 cm 2394 What is the speed of sound in oxygen? 2395 GPT-40: 2397 To determine the speed of sound in oxygen using the given setup, we will analyze the standing wave 2398 pattern formed in the Kundt's tube. Here's the step-by-step reasoning: 2399 2400 Step 1: Understand the standing wave pattern 2401 2402 Step 2: Relationship between the wavelength and the distance between nodes 2403 2404 Step 3: Determine the wavelength 2405 - From the diagram, the total length of the tube is \(123 \, \text{cm} = 1.23 \, \text{m}\). - The cork particles form regularly spaced piles, indicating the positions of the nodes. 2406 Count the number of node-to-node segments in the tube. From the diagram, there are 5 segments. 2407 - The total length of the tube corresponds to 5 half-wavelengths:  $[5 \cdot f^{12} = 1.23 , \text{m}]$ Solve for \(\lambda\) (the wavelength): 2409  $[\lambda = \frac{2 \cdot 1.23}{5} = 0.492 \, \text{text}]$ 2410 Step 4: Use the wave equation to find the speed of sound 2411 2412 Final Answer: 2413 The speed of sound in oxygen is approximately: \[\boxed{196.8 \, \text{m/s}}\] 2414 2415 Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning 2416 2417 **Error Type: Visual Reasoning Errors** 2418 2419 Ground Truth: 1.2 m/s 2420 2421

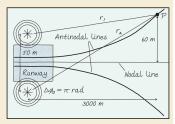
Figure 18: A sample error case of Wave/Acoustics. Error category: Visual Reasoning Error Back to List of Figures | Back to Table Index

# Wave/Acoustics



#### Description:

Your team is designing a landing aid for pilots in rain or fog by placing two radio transmitters 50 m apart on either side of a runway. These transmitters emit the same frequency but are out of phase, creating a nodal line down the runway center. The pilot hears nothing when aligned, but hears a beep if off-center. For accurate guidance, the first intensity maxima should appear 60 m from the nodal line at a distance of 3.0 km.



#### Question:

What frequency should you specify for the transmitters?

#### GPT-40:

To determine the frequency of the transmitters, we need to analyze the interference pattern created by the two sources. The goal is to ensure that the first antinodal line (intensity maxima) is  $\langle (60 \, \backslash \, \text{mathrm{m}} \, \rangle )$  away from the nodal line at a distance of  $\langle (3.0 \, \backslash \, \text{mathrm{km}} \, \rangle )$  from the transmitters.

#### Step 1: Geometry of the problem

At a distance \( L = 3000 \, \mathrm{m} \) from the transmitters, the first antinodal line is \( y = 60 \, \mathrm{m} \) away from the nodal line. This forms a right triangle with the base \( L \) and height \( y \).

Step 2: Path difference for the first antinodal line

#### Step 3: Wavelength of the radio waves

For the first antinodal line, the path difference  $\setminus ( \cdot \mid r \mid )$  corresponds to one wavelength:

 $[ \beta r = \beta . ]$ 

Thus: \[ \lambda = 1.0 \, \mathrm{m}. \]

#### Step 4: Frequency of the radio waves

The frequency  $\ \ f \ \$  of the radio waves is related to the wavelength  $\ \ \ \$  and the speed of light  $\ \ \ \$  by:  $\ \ f = \frac{c}{\lambda}.$ 

Using  $(c = 3.0 \times 10^8 , \mathbf{m}_{m/s})$  and  $(\lambda = 1.0 , \mathbf{m}_{m})$ :

 $[f = \frac{3.0 \times 10^8}{1.0} = 3.0 \times 10^8 , \mathbb{Hz}. ]$ 

#### Final Answer:

The frequency of the transmitters should be:  $\[ \boxed{300 \, \mathrm{MHz}} \]$ 

#### Reasoning Type: Spatial Relation Reasoning, Multi-Formula Reasoning

**Error Type: Text Reasoning Errors** 

Ground Truth: 150 MHz

Figure 19: A sample error case of Wave/Acoustics. Error category: Text Reasoning Error Back to List of Figures | Back to Table Index

#### Wave/Acoustics Description: ₹3.00 m Two loudspeakers are placed 3.00 m apart, as shown in figure. They emit 474 Hz sounds, in phase. A microphone is placed 3.20 m distant from a point 3.20 m midway between the two speakers, where an intensity maximum is recorded. \\ Question: How far must the microphone be moved to the right to find the first intensity GPT-40: To solve this problem, we need to determine the distance the microphone must be moved to the right to encounter the first intensity minimum. This involves understanding the conditions for destructive interference of sound waves. Step 1: Conditions for Destructive Interference Step 2: Wavelength of the Sound Step 3: Geometry of the Problem Step 4: Step 4: Approximation for Small \( x \) For small $\ (x \ )$ , we can use a first-order approximation for the square root: $\ [ \sqrt{a^2 + b^2} \ approx \ a + \frac{b^2}{2a}, \quad \ text[if ] \ b \ ll \ a. \ ]$ $+ \frac{(1.50 + x)^2}{2(3.20)},$ Step 5: Solve for $\setminus (x \setminus)$ Final Answer: The speed of the waves in the coffee is approximately: $[\boxed{0.08 \, \mathrm{m/s}}\]$ Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning Error Type: Lack of Knowledge Ground Truth: 0.429m

Figure 20: A sample error case of Wave/Acoustics. Error category: Lack of Knowledge Back to List of Figures | Back to Table Index

2538 2539 2540 2541 2542 2543 **Optics** 2544 2545 2546 Description: 2547  $heta_1$  ; A ray of light travels from air into another medium, 2548 making an angle of  $\theta$ 1=45.0° with the normal as in 2549 Air 2550 Question: 2551 Second Find the angle of refraction  $\theta 2$ , if the second 2552 medium medium is water.  $\theta_2$ 2553 2554 2555 GPT-40: 2557 To solve this problem, we will use \*\*Snell's Law\*\*, which relates the angles of incidence and refraction to the refractive indices of the two media:  $[n_1 \sin \theta_1 = n_2 \sin \theta_2]$ 2559  $- (n_1)$  is the refractive index of the first medium (air), 2560  $- (n_2)$  is the refractive index of the second medium (water), 2561 - \(\theta\_1\) is the angle of incidence, 2562 2563 Step 1: Assign known values 2564 - The refractive index of air is approximately  $(n_1 = 1.00)$ , - The refractive index of water is  $(n_2 = 1.33)$ , 2565 - The angle of incidence is  $\ ( \text{theta}_1 = 45.0 \ \ ).$ 2566 2567 Step 2: Apply Snell's Law Substitute the known values into Snell's Law:  $[1.00 \cdot \sin(45.0^\circ) = 1.33 \cdot \sinh(\frac{2})]$ 2568 2569 Step 3: Simplify the equation The sine of \( 45.0^\circ \) is: \[ \sin(45.0^\circ) = \frac{\sqrt{2}}{2} \approx 0.707 \] 2570 Substitute this value into the equation:  $[1.00 \cdot 0.707 = 1.33 \cdot (theta_2)] [0.707 = 1.33 \cdot (theta_2)]$ 2571 \cdot \sin(\theta\_2) \] 2572 Step 4: Solve for \(\\sin(\theta\_2)\) 2573 Divide both sides of the equation by  $(1.33 \): \[ \cdot \sin(\theta_2) = \frac{0.707}{1.33} \] \[ \cdot \sin(\theta_2) = \frac{0.707}{1.33} \]$ 2574 \approx 0.531 \] 2575 Step 5: Find \(\theta\_2\) 2576 Take the inverse sine (\(\arcsin \)) of \( 0.531 \) to find \(\theta\_2 \): \[ \theta\_2 = \arcsin(0.531) \] 2577 Using a calculator: \[ \theta\_2 \approx 32.0^\circ \] 2578 Final Answer: 2579 The angle of refraction is:  $\[ \boxed{\hat _2 \approx 32.0^\circ \] }$ 2580 2581 Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning 2582

Figure 21: A sample correct case of Optics. Back to List of Figures | Back to Table Index

2583

2584 2585

2586

258825892590

Ground Truth: 32.0°

#### **Optics** Description: Ouestion: Assume a transparent rod with diameter d=2.00 Determine the maximum angle $\theta$ at which light rays µm and an index of refraction of 1.36. can be incident on the end of the rod (as shown in the figure) such that they undergo total internal reflection along the walls of the rod. Options: A. 72.2° B. 65.4° C. 67.2° D. 60.0° dGPT-40: To determine the maximum angle \(\\\\), we use the concept of total internal reflection. The critical where $(n_1 = 1.36)$ (index of refraction of the rod) and $(n_2 = 1.00)$ (index of refraction of air). 47.8^\circ \] The maximum angle \(\\\\) for total internal reflection is related to the critical angle by: \[ \theta = 90^\circ - \theta\_c \] \[ \theta = 90^\circ - 47.8^\circ = 42.2^\circ \] However, this is the angle relative to the normal. For the angle relative to the axis of the rod, we use: \[ \theta = 90^\circ - 42.2^\circ = 67.2^\circ \] Thus, the correct answer is: C Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning Ground Truth: C. 67.2

Figure 22: A sample correct case of Optics. Back to List of Figures | Back to Table Index

**Optics** Description: When the light ray illustrated in figure passes through the glass block of index of refraction n = 1.50, it is shifted laterally by the distance d."  $\overline{\theta}$ Question: Find the time interval required for the light to pass through the glass 45.0° block. GPT-40: To solve this problem, we will use Snell's Law and the geometry of the prism. The steps are as follows: Step 1: Understand the geometry of the prism and the light path Step 2: Relate the angles inside the prism Inside the prism, the beam travels at an angle of \( 45^\circ \) with respect to the normal to the hypotenuse. This is because the prism is isosceles, and the beam enters perpendicular to one leg. - At the hypotenuse, the beam refracts and exits the prism. Let the angle of refraction (outside the prism) be  $\ \ (r \ )$ , measured with respect to the normal to the hypotenuse. From the geometry of the problem: \[\theta = 45^\circ - r\] Substitute \(\theta = 15.0^\circ \): \[15.0^\circ = 45^\circ - r\]  $[r = 30.0^circ]$ Step 3: Apply Snell's Law at the hypotenuse Step 4: Solve for \( n \) Final Answer: The index of refraction of the prism is: \[\boxed{1.414}\] Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning **Error Type: Visual Reasoning Errors** Ground Truth: 1.2m/s 

Figure 23: A sample error case of Optics. Error category: Visual Reasoning Error Back to List of Figures | Back to Table Index

**Optics** 30.0° Description: When the light ray illustrated in figure passes through the glass block of index of refraction n = 1.50, it is shifted laterally by the distance d. 2.00 cm Question: Find the time interval required for the light to pass through the glass block. GPT-40: To find the time interval required for the light to pass through the glass block, we need to calculate the time it takes for the light to travel the distance it covers inside the block. Here's the step-by-step reasoning: Step 1: Understand the problem Step 2: Speed of light in the glass Step 3: Distance traveled by the light in the glass The light ray travels diagonally through the glass block. The thickness of the block is \( (2.00 \, \text{cm} = 0.0200 \, \text{m} \), and the light travels at an angle of \( 30^\circ \) with respect to the normal. The actual distance \(\((L\)\)\) traveled by the light inside the block is the hypotenuse of a right triangle, where the thickness of the block is the adjacent side. Using trigonometry:  $[L = \frac{\text{text{thickness}}{\cos(30^\circ)}.]}$ Substitute the values:  $[L = \frac{0.0200}{\cos(30^\circ circ)}.]$ Step 4: Time taken to travel through the glass Final Answer: The time interval required for the light to pass through the glass block is approximately: \[t \approx 1.16 \times 10^{-10} \, \text{s}.\] Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning **Error Type: Text Reasoning Error** Ground Truth: 106ps 

Figure 24: A sample error case of Optics. Error category: Text Reasoning Error Back to List of Figures | Back to Table Index

**Optics** 12.0 cm Description: A narrow beam of ultrasonic waves reflects off the liver tumor illustrated in 50.0° figure. The speed of the wave is 10.0% less in the liver than in the surrounding medium. Question: Liver Determine the depth of the tumor. Tumor GPT-40: To determine the depth of the tumor, we will use the information provided in the diagram and the relationship between the speed of the ultrasonic wave in the liver and the surrounding medium. Here's the step-by-step reasoning: Step 1: Understand the geometry of the problem Step 2: Relate the geometry to the depth The wave travels in a triangular path. The horizontal distance between the entry and exit points is \( 12.0 \, \text{cm} \), and the wave reflects symmetrically. This means the horizontal distance from the entry point to the reflection point is \( 6.0 \, \text{cm} \). Using trigonometry, the depth  $\setminus$  (d  $\setminus$ ) of the tumor can be related to the horizontal distance and the angle of incidence:\[\tan(50.0^\circ) = \frac{d}{6.0 \, \text{cm}}.\] Rearranging for  $\ (\ d\ )$ :  $\ [\ d = 6.0\ , \ text{cm} \ \ tan(50.0^\circ).\ ]$ Step 3: Step 3: Calculate \( d \) Step 4: Adjust for the speed difference Final Answer: The depth of the tumor is approximately: \[\boxed{7.94 \, \text{cm}}.\] Reasoning Type: Spatial Relation Reasoning, Physical Model Grounding Reasoning Error Type: Lack of Knowledge Ground Truth: 6.30cm 

Figure 25: A sample error case of Optics. Error category: Lack of Knowledge Back to List of Figures | Back to Table Index

#### 2811 2812 2813

### 2814 2815

#### 2816 2817 2818

# 2819

#### 2821 2822

## 2827

2829

2831

2833

2836

2837

2839

2841 2842

2843

2845 2846 2847

2848 2849 2850

2851 2852

2856 2857

2859

# Modern Physics



300 m

 $(x_{R}, t_{R}) = (300 \text{ m}, 0 \text{ s})$ 

#### Description:

Peggy is standing in the center of a long, flat railroad car that has firecrackers tied to both ends. The car moves past Ryan, who is standing on the ground, with velocity v = 0.8c. Flashes from the exploding firecrackers reach him simultaneously  $1.0\mu s$  after the instant that Peggy passes him, and he later finds burn marks on the track 300m to either side of where he had been standing.

#### **Ouestion:** According to Peggy, what times do the explosions occur relative

# Event L -300 m

 $(x_1, t_1) = (-300 \text{ m}, 0 \text{ s})$ 

Peggy

#### GPT-40:

To solve this problem, we will use the principles of special relativity, specifically the Lorentz transformation equations, to determine the times of the explosions in Peggy's frame ((S')) relative to the time that Ryan passes her ((t' = 0)).

Step 1: Understand the problem and given data

Rvan's frame  $(\(S\))$ :

to the time that Ryan passes her?

- The firecrackers explode at  $(x_L = -300\;m\setminus)$  and  $(x_R = 300\;m\setminus)$ , both at  $(t = 0\setminus)$ .
- The flashes from the explosions reach Ryan simultaneously at (t = 1.0); mu s.
- The velocity of Peggy's frame ((S')) relative to Ryan's frame ((S)) is (v = 0.8c). Peggy's frame (\(S'\)):
- Peggy is at rest in her frame, and the firecrackers are tied to the ends of the railroad car.
- We need to find the times of the explosions ( $(t'_L)$ ) and  $(t'_R)$ ) in Peggy's frame.

#### Step 2: Lorentz transformation equations

The Lorentz transformation equations relate the coordinates  $\backslash ((x, t)\backslash)$  in Ryan's frame  $(\backslash (S\backslash))$  to the coordinates ((x', t')) in Peggy's frame ((S')):  $[t' = \gamma \left(t - \frac{vx}{c^2} \right), ] [x' = \gamma \left(t - \frac{vx}{c^2} \right)$  $\gamma \left( x - vt \right),$ 

where:  $\[ \gamma = \frac{1}{\sqrt{1 - v^2/c^2}}. \]$ 

For  $\langle v = 0.8c \rangle$ , we calculate:  $[\gamma = \frac{1}{\sqrt{1 - (0.8)^2}} = \frac{1}$  ${\sqrt{0.36}} = \frac{5}{3}.$ 

Step 3: Apply the Lorentz transformation for each event

Event L (explosion at  $\langle x_L = -300 \rangle$ ,  $\langle t_L = 0 \rangle$ ):

 $\label{torsequence} Using the Lorentz transformation for time: $$ \{t'_L = \gamma \left(t_L - \frac{v_L}{c^2} \right). $$$ Substitute the values:  $[t'_L = \frac{5}{3} \left(0.8c)(-300\mbox{ })\c^2\right)$ 

Since  $(c = 3 \times 10^8);m/s)$ , we calculate:  $[t'_L = \frac{5}{3} \cdot 0 \times 10^8] = 10^8$ Event R (explosion at  $(x_R = 300; m)$ ),  $(t_R = 0)$ :

Using the Lorentz transformation for time:  $\[t'_R = \gamma \]$  \[ \text{t\_R - \frac{v x\_R}(c^2} \right] \] Substitute the values:  $[t'_R = \frac{5}{3} \left(0 - \frac{(0.8c)(300\mb)}{c^2} \right).$ 

 $10^{8} = -\frac{10^{8}} = -\frac{10^{$ -1.33\;\mu s. \]

#### Step 4: Interpret the results

- In Peggy's frame (\(S'\)):
- The explosion at  $(x'_L)$  (left firecracker) occurs at  $(t'_L = 1.33)$ ; mu s\).
- The explosion at  $(x'_R)$  (right firecracker) occurs at  $(t'_R = -1.33)$ ; mu s).

This means that, according to Peggy, the right firecracker explodes before Ryan passes her, and the left firecracker explodes after Ryan passes her.

Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning

Ground Truth: -1.33 us

Figure 26: A sample correct case of Modern Physics. Back to List of Figures | Back to Table Index

# **Modern Physics**



#### Description:

# An enemy spaceship is moving toward your starfighter with a speed, as measured in your frame, of 0.400c. The enemy ship fires a missile toward you at a speed of 0.700c relative to the enemy ship. You measure that the enemy ship is $8.00 \times 10^6$ km away from you when the missile is fixed

#### Question:

If you measure that the enemy ship is  $8.00\times10^6 km$  away from you when the missile is fired, how much time, measured in your frame, will it take the missile to reach you?

#### Options:

A. 31.0s

B. 30.4s

C. 33.1s

D. 32.8s



#### GPT-40:

To solve this, we use the relativistic velocity addition formula to find the speed of the missile relative to the starfighter:

 $[v_{\text{missile/enemy}}] = 0.700c, \quad v_{\text{enemy/starfighter}} = 0.400c ]$ 

 $\label{eq:conditional_condition} $$ v_{\text{missile/starfighter}} = \frac{0.700c + 0.400c}{1 + \frac{(0.700c)(0.400c)}{c^2}} \] $$$ 

 $\label{eq:v_starfighter} $$ \left[ v_{\text{missile/starfighter}} = \frac{1.100c}{1 + 0.280} = \frac{1.100c}{1.280} \right] $$$ 

Now, calculate the time it takes for the missile to travel the distance of  $(8.00 \times 10^6 , \text{km})$  (or  $(8.00 \times 10^9 , \text{km})$ ) at this speed:

Substitute  $(c = 3.00 \times 10^8 , \text{m/s})$ :

 $[t \rightarrow 10^1 \rightarrow 10$ 

The correct answer is: A

Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning

Ground Truth: A. 31.0s

Figure 27: A sample correct case of Modern Physics.

Back to List of Figures | Back to Table Index

#### 2916 2917 2918 **ModernPhysics** 2919 2920 2921 Description: 2922 An experiment was performed in which neutrons were shot through two slits spaced 0.10 nm apart and detected 3.5m 2924 behind the slits. Figure shows the detector output. Neutron intensity 2925 2926 2927 Question: 2928 To one significant figure, what was the speed of the 2929 neutrons? 2930 2931 GPT-40: To determine the speed of the neutrons, we can use the principles of wave interference and de Broglie's 2933 hypothesis. Here's the step-by-step reasoning: 2934 Step 1: Understand the setup 2935 10^{-10} \, \text{m} \) apart. The interference pattern is detected on a screen located \( L = 3.5 \, \text{m} 2937 \) behind the slits. The distance between adjacent maxima in the interference pattern is \(\\Delta y = 100\) 2939 Step 2: Use the double-slit interference formula The position of the maxima in a double-slit interference pattern is given by: 2940 $[y_m = m \frac{\label{y_m}}{\label{y_m}}]$ 2941 2942 Step 3: Relate the wavelength to the speed of the neutrons 2943 Using de Broglie's hypothesis, the wavelength of a particle is related to its momentum by: 2944 \[ \lambda = \frac{h}{p}, \] where: 2945 - \( h = $6.63 \times 10^{-34} \, \text{J-s} \)$ is Planck's constant, 2946 - \( p = mv \) is the momentum of the neutron, 2947 - $\ \ \ = 1.675 \times 10^{-27} \ , \ \ \ )$ is the mass of a neutron, - $\langle (v \rangle)$ is the speed of the neutron. 2948 2949 Rearranging for $\langle (v \rangle)$ : 2950 $[v = \frac{h}{m \lambda}.]$ 2951 Substitute the known values: 2952 \, \text{m})}. \] 2953 $[v = 1.38 \times 10^3 , \text{x}_{m/s}.]$ 2954 2955 Step 4: Round to one significant figure To one significant figure, the speed of the neutrons is: 2956 $[v \cdot 1 \cdot 1^3 \cdot 1^3 \cdot v \cdot 1^3 \cdot 1^3$ 2957 Final Answer: 2958 The speed of the neutrons is approximately $\ (\mathbf{1} \times 10^3 \, , \text{x}).$ 2959 2960 Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning 2961 2962 **Error Type: Visual Reasoning Error** 2963 2964 Ground Truth: 200m/S 2965

Figure 28: A sample error case of Modern Physics. Error category: Visual Reasoning Error Back to List of Figures | Back to Table Index

2967

#### 2970 2971 2972 2973 **ModernPhysics** 2974 2975 2976 Description: 2977 $1.80 \times 10^{12}$ r Owen and Dina are at rest in frame S, which is moving with 2978 0.800crespect to frame S. They play a game of catch while Ed, at 2979 rest in frame S, watches the action as shown in figure. Owen throws the ball to Dina. 2980 Dina 2981 2982 Ed 2983 Question: 2984 What time interval is required for the ball to reach Dina? 2985 GPT-40: 2987 To determine the time interval required for the ball to reach Dina, we need to analyze the situation using the principles of special relativity. Here's the step-by-step reasoning: 2989 Step 1: Understand the problem and given data\*\* Frames of reference:\* 2991 - Frame \( S' \): The rest frame of Owen and Dina. - Frame \( S \): The rest frame of Ed. 2992 2993 Given data:\*\* 2994 The relative velocity of frame $\ (S')$ with respect to frame $\ (S): \ (v = 0.600c)$ . - The velocity of the ball relative to (S'): (u' = 0.800c) (from Owen to Dina). 2995 - The distance between Owen and Dina in $\ (S' \): \ (d' = 1.80 \times 10^{12} \, \text{text}_m) \).$ 2996 We are tasked with finding the time interval for the ball to reach Dina \*\*as observed in frame \( S \)\*\*. 2997 2998 Step 2: Relate the velocity of the ball in $\backslash ($ S $\backslash )$ to the velocity in $\backslash ($ S' $\backslash )^{**}$ 2999 The velocity of the ball in frame $(S \setminus (u \setminus u))$ can be calculated using the relativistic velocity addition formula: 3000 3001 $[ u = \frac{u' + v}{1 + \frac{u'v}{c^2}} ]$ 3002 Substitute the given values: 3003 - \( u' = 0.800c \) 3004 - \( v = 0.600c \) 3005 3006 3007 Step 3: Calculate the time interval in frame \( S \) 3008 3009 Final Answer: 3010 The time interval required for the ball to reach Dina, as observed in frame \( S \), is approximately: 3011 $\[ \boxed{6.34 \, \text{seconds}} \]$ 3012 3013 Reasoning Type: Physical Model Grounding Reasoning, Multi-Formula Reasoning 3014 3015 **Error Type: Text Reasoning Errors** 3016 3017 Ground Truth: 4.88 \* 10^3 s 3018

Figure 29: A sample error case of Modern Physics. Error category: Text Reasoning Error Back to List of Figures | Back to Table Index

3019

3020

#### **ModernPhysics** Description: Figure shows the wave function of an electron in a rigid $\psi(x)$ box. The electron energy is 6.0eV. Question: How long is the box? To determine the length of the box, we will use the quantum mechanical model of a particle in a rigid box. The energy levels of an electron in a rigid box are given by: $\[ E_n = \frac{n^2 h^2}{8mL^2} \]$ where: - \( $E_n$ \) is the energy of the electron in the \( n \)-th energy level, - $\ (m \ )$ is the mass of the electron ( $\ (m = 9.109 \ ), \ (31) \ , \ (kg) \ )),$ - $\setminus$ ( L $\setminus$ ) is the length of the box. Step 1: Identify the quantum number $\setminus$ ( n $\setminus$ ) Step 4: Convert to nanometers $[ L = 7.52 \times 10^{-10} \, \text{m} = 0.752 \, \text{m}. ]$ Final Answer: The length of the box is approximately: \[ \boxed{0.752 \, \text{nm}} \] Reasoning Type: Physical Model Grounding Reasoning, Spatial Relation Reasoning Error Type: Lack of Knowledge Ground Truth: 1 nm

Figure 30: A sample error case of Modern Physics. Error category: Lack of Knowledge Back to List of Figures | Back to Table Index

DATA ANNOTATION PROTOCOL

3078

3079 3080

3082

3084 3085

3086 3087

3088 3089

3090

3091

3092

3093

visual context.

**Types of Questions:** 

DATA COLLECTION

I.1

#### 3096 **Image Types.** The annotators should find images with realistic physical senarios. 3097 I.2 GENERAL GUIDELINES 3099 • General Principles: Annotations should be accurate and uniform, and maintain a high level 3100 of academic quality. 3101 • Specific Instructions: 3102 3103 - All questions should be written in English. 3104 All questions must contain one physical image. 3105 - All images in question should be realistic, in specific physical scenarios. 3106 - The question should not be ambiguous and can be answered with one of the given 3107 options or a short answer. 3108 - Annotate all data fields, including the description, simplified description, question, 3109 answer options, the correct answer, image, and domain. 3110 3111 I.3 DATA FORMAT AND STRUCTURE 3112 • **JSON File Format:** The structured JSON format will include fields for index number, 3113 description, simplified description, question, answer options, correct answer, and domain. 3114 3115 Naming Conventions: 3116 Each collected sample will be stored on a single line in a JSONL file. 3117 - Image files should follow a standard naming rule: {QuesNum}.png 3118 • Interleaving Question with Images: The images should be inserted as a file path in the 3119 question. 3120 3121 I.4 QUALITY CONTROL AND VALIDATION 3122 3123 Annotators will cross-check each other's work to ensure accuracy and compliance with the annotation guidelines. 3124 3125 Periodic reviews of randomly selected samples from the dataset will be carried out to 3126 maintain consistent quality over time. 3127 3128 I.5 HANDLING AMBIGUITIES 3129 Any ambiguous or unclear data entries should be marked for thorough review. Such questions will 3130 be collectively discussed during team meetings to develop a consistent and standardized annotation 3131 strategy. 58

This document outlines a detailed procedure for annotating a dataset of physics questions that include

**Sources of Data.** Data is collected from freely accessible online resources, textbooks, and other materials. Annotators are instructed to use a wide range of sources rather than relying on just one.

required to create a corresponding open-ended version of the same problem.

• Multiple-Choice Questions: These consist of a question accompanied by four answer

options, with only one being correct. For each multiple-choice question, annotators are also

• Open-Ended Questions: These include formats such as short-answer and calculation-

based problems. Questions with excessively lengthy answers should be avoided. For each open-ended question, a corresponding multiple-choice version should also be constructed.

#### I.6 ETHICAL CONSIDERATIONS

- Copyright and Licensing: Annotators must strictly follow all applicable copyright and licensing rules. Content from sources that restrict reproduction or redistribution will be excluded without exception.
- **Data Privacy:** Upholding data privacy and ethical standards is essential. Annotators should refrain from including any questions that involve personal or sensitive information.

#### I.7 DATA CONTAMINATION CONSIDERATIONS

When developing benchmarks for evaluating foundation models, it is crucial to account for the potential risk of data contamination. To mitigate this, annotators should deliberately avoid simple questions with widely available answers. Instead, they should prioritize selecting problems whose solutions are embedded in less conspicuous places—such as in supplementary materials or at the end of lengthy textbooks. This strategy helps ensure that the benchmark effectively challenges models to demonstrate genuine comprehension and reasoning across complex and less accessible content.