# LOCA-R: Near-Perfect Performance on the Chinese Physics Olympiad 2025

**Dong-Shan Jian**[1,*] **Xiang Li**[1,*] **Chen-Xu Yan**[1] **Hui-Wen Zheng**[1]
**Zhi-Zhang Bian**[1] **You-Le Fang**[1] **Sheng-Qi Zhang**[1,2]
**Bing-Rui Gong**[1] **Ren-Xi He**[1] **Jing-Tian Zhang**[1] **Ce Meng**[1]
**Yan-Qing Ma**[1,2,†]
[1]School of Physics, Peking University, Beijing 100871, China
[2]Center for High Energy Physics, Peking University, Beijing 100871, China
{dsjian}@stu.pku.edu.cn
{lix-PHY,yqma}@pku.edu.cn

## Abstract

Olympiad-level physics problem-solving presents a significant challenge for both humans and artificial intelligence (AI), as it requires a sophisticated integration of precise calculation, abstract reasoning, and a fundamental grasp of physical principles. The Chinese Physics Olympiad (CPhO), renowned for its complexity and depth, serves as an ideal and rigorous testbed for these advanced capabilities. In this paper, we introduce LOCA-R (LOgical Chain Augmentation for Reasoning), an improved version of the LOCA framework adapted for complex reasoning, and apply it to the CPhO 2025 theory examination. LOCA-R achieves a near-perfect score of 313 out of 320 points, solidly surpassing the highest-scoring human competitor and significantly outperforming all baseline methods.[1]

## 1 Introduction

Physics problem-solving, or more generally complex scientific reasoning, stands as a challenging frontier for both humans and AI. It demands not only mathematical derivations, but also the ability to translate complex real-world scenarios described in natural language into abstract models. This process requires a deep understanding of physical laws to select and apply the appropriate principles. While Large Language Models (LLMs) have demonstrated remarkable success in general domains or structured fields like mathematics and coding(Brown et al., 2020; OpenAI, 2022; Achiam et al., 2023; Anil et al., 2023; Touvron et al., 2023a;b; Liu et al., 2024a; Guo et al., 2025; Anthropic, 2025; Comanici et al., 2025; OpenAI, 2025a;b; Team, 2025; Team et al., 2025; Yang et al., 2025), their performance in physics has fallen behind as shown in many prior works(He et al., 2024; Chow et al., 2025; Chung et al., 2025a; Feng et al., 2025; Qiu et al., 2025b; Zhang et al., 2025c;b; Xu et al., 2025; Siddique et al., 2025; Carrit Delgado Pinheiro et al., 2025). Therefore, given the strict demand for accuracy in scientific problems, the most critical task is to maximize the problem-solving capabilities of LLMs, pushing them towards perfect scores. The trade-off between performance and computational cost, while important, is a subsequent consideration.

To rigorously assess and push the boundaries of AI's capabilities in this challenging domain, a suitable benchmark is essential. The Chinese Physics Olympiad (CPhO) (CPS, 2025) is a premier national

---

[*]These authors contributed equally to this work.

[†]Corresponding author.

[1]Github repository: https://github.com/Science-Discovery/LOCA-R.

physics competition organized annually in China. Much like the prestigious International Physics Olympiad (IPhO)(IPhO, 2025), CPhO is designed to challenge the country's senior high students in physics. The 42nd CPhO in 2025, with its final round held in Fuzhou at the end of October, represents the highest level of pre-university physics education and problem-solving in the nation. For our evaluation, we focus on the theory examination of CPhO 2025 for two key reasons. First, CPhO problems are known for their depth, often requiring contestants to synthesize knowledge from different areas of physics and apply creative thinking to unfamiliar scenarios. Second, unlike many existing datasets that may suffer from data contamination issues, CPhO 2025 problems are newly created, ensuring a fair and novel test.

State-of-the-art LLMs possess sufficient calculation abilities and foundational knowledge for complex physics(Qiu et al., 2025b; Chung et al., 2025a), yet their performance remains limited. This indicates that the primary bottleneck is not a lack of information, but rather the failure of generic reasoning-enhancement strategies in this domain. A more effective approach, therefore, is to cultivate the model's intrinsic physics reasoning rather than relying on one-size-fits-all methods(Chung et al., 2025b; Zhang et al., 2025a; Wang et al., 2024; Xi et al., 2025).

To address this issue, we introduce LOCA-R (LOgical Chain Augmentation for Reasoning), a framework that structures and refines LLM reasoning in an atomic manner. LOCA-R builds upon the recently proposed LOCA framework (Fang et al., 2025). The core idea of LOCA, originally designed for filtering physics benchmarks, is to enforce the verifiability of reasoning by decomposing each step into a foundational principle and its subsequent derivation. To tackle complex Olympiad-level problems, we enhance this framework's robustness by introducing a more rigorous, hard-coded atomic and sequential review mechanism and a dedicated problem interpretation module. These improvements ensure greater logical coherence and a deeper initial understanding of the problem.

When applied to the CPhO 2025 theory examination, LOCA-R achieves a near-perfect score of 313 / 320. This result is unprecedented: LOCA-R's score on the theory section alone surpasses the total score (theory + experimental) of the best-performing human gold medalist (273 / 400). Furthermore, its performance significantly exceeds that of all baseline LLM reasoning methods. The solutions generated by LOCA-R are not only correct but also highly readable, presenting clear, step-by-step reasoning that holds great promise for educational applications and the training of future models.

## 2 RELATED WORKS

**LLM Reasoning.** LLMs have been broadly applied to address general reasoning tasks by generating intermediate steps and exploring various reasoning paths, and are sometimes equipped with external tools such as a code interpreter. Existing strategies can empirically be divided into three main categories:

- Test-time scaling frameworks, like Chain-of-Thought (CoT)(Kojima et al., 2022; Wei et al., 2022), Tree-of-Thought (ToT) (Yao et al., 2023), Graph-of-Thought (GoT) (Besta et al., 2024) and self-refine(Madaan et al., 2023; Liu et al., 2024b; Liang et al., 2024; Zhang et al., 2024), usually relying only on the model's internal capabilities.

- Agentic frameworks (usually with external tools), such as Plan-and-Solve(Wang et al., 2023), Re-Act(Yao et al., 2022) and Multi-Agent Debate (MAD)(Du et al., 2023).

- Training-based methods like Supervised Fine-Tuning (SFT) or Reinforcement Learning (RL)(Zhang et al., 2024; Lu et al., 2022; Lewkowycz et al., 2022), which critically depend on the availability of large-scale, expert-validated training data and can be prohibitively costly.

Targeted approaches like Physics Supernova(Qiu et al., 2025a) , Physics Reasoner(Pang et al., 2025) and Mixture of Refinement Agents (MoRA)(Jaiswal et al., 2024) have been developed for physics. Despite

this, the fundamental issues of modeling inaccuracies and the selection of relevant principles remain largely unaddressed in the literature.

**Physics Benchmarks.** The effort to advance AI in physics problem-solving is broadly supported by specialized benchmarks, ranging from expert-curated datasets focusing on competition-level difficulty, such as PHYBench (Qiu et al., 2025b) and OlympiadBench (He et al., 2024), to large-scale collections reflecting university-level coursework, such as PHYSICS (Feng et al., 2025), ABench-Physics(Zhang et al., 2025c), PhysReason(Zhang et al., 2025b), and PhysicsEval(Siddique et al., 2025), and to specific domains such as TPBench(Chung et al., 2025a). However, the CPhO 2025 offers distinct advantages, including the absence of data contamination and the availability of detailed, reliable reference answers.

## 3 METHOD

The methodology builds upon the recently proposed LOCA framework (Fang et al., 2025). Originally designed for filtering physics benchmarks, LOCA's core idea of structured augmentation and review also proves highly effective for complex problem-solving tasks. Through targeted adaptations and improvements, we develop LOCA-R, a powerful LLM-based reasoning framework for tackling Olympiad-level physics problems. In this section, we first provide a brief overview of the foundational concepts of LOCA. We then detail the crucial improvements we introduce: a hard-coded atomic and sequential review mechanism, and a dedicated problem interpretation module.

### 3.1 PRELIMINARY: A REVIEW ON LOGICAL CHAIN AUGMENTATION

The cornerstone of LOCA is logical chain augmentation, a mechanism that transforms a raw, unstructured solution into a detailed and structured logical chain. This process addresses two common weaknesses: **non-atomicity**, where a single step combines multiple reasoning acts into a "logical leap", and **implicit justification**, where the underlying scientific principle is mixed with its subsequent mathematical derivation or omitted entirely. To resolve these issues, LOCA remaps a raw solution $S_{\text{raw}}$, represented as a sequence of steps $(s_1, \ldots, s_n)$, to a structured, augmented solution $S_{\text{aug}}$ through two key operations: chain completion and structured decomposition.

**Chain Completion.** This operation enforces atomicity by interpolating missing logical steps. A step $s_i : C_{i-1} \to C_i$ (where $C$ is the context) is considered non-atomic if it implicitly contains an intermediate context $C_{\text{int}}$. LOCA decomposes each such step $s_i$ into a more fundamental subsequence $(s'_{i,1}, \ldots, s'_{i,k})$:

$$\forall s_i \in S_{\text{raw}}, \text{ if } \neg\text{IsAtomic}(s_i), \text{ then } s_i \mapsto (s'_{i,1}, \ldots, s'_{i,k}). \tag{1}$$

This results in a new, more detailed sequence (solution) $S' = (s'_1, \ldots, s'_m)$ with $m \geq n$.

**Structured Decomposition.** To enhance clarity and verifiability, LOCA decomposes each atomic step $s'_j$ into two distinct components, forming a structured tuple $s'_j = (P_j, D_j)$. The **principle** ($P_j$) defines the logical foundation of the step, such as a physical law or a mathematical theorem. The **derivation** ($D_j$) then describes the specific application of this principle to produce the step's outcome.

The final augmented solution is thus a sequence of these structured tuples:

$$S_{\text{aug}} = ((P_1, D_1), \ldots, (P_m, D_m)). \tag{2}$$

This structured representation enables a highly reliable augment-and-review loop. The correctness of the solution is assessed by two specialized review agents—one for the principle ($\mathcal{R}_P$) and one for the derivation

$(\mathcal{R}_D)$—with the final judgment integrating the feedback of both agents:

$$\mathcal{R} = \mathcal{R}_P \wedge \mathcal{R}_D. \tag{3}$$

This decomposition is crucial for physics reasoning as it isolates distinct sources of error.
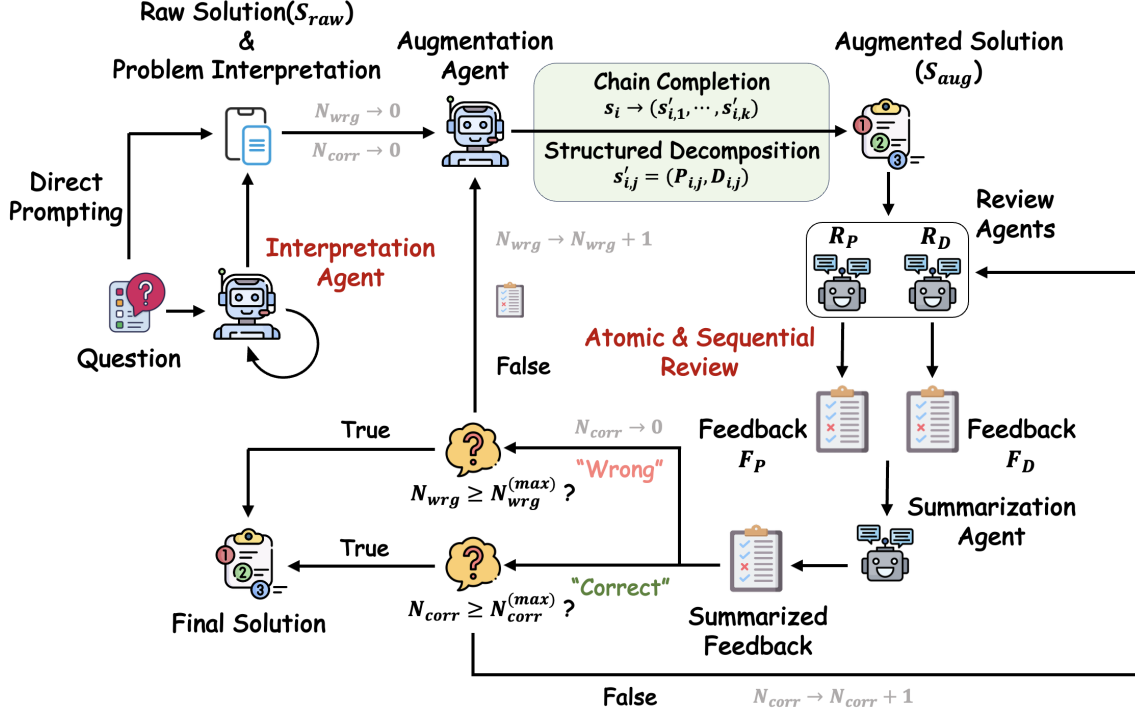


Figure 1: **An overview of LOCA-R's architecture.** The framework builds upon logical chain augmentation by implementing an iterative augment-and-review loop. It is further enhanced with an atomic, sequential review mechanism and a dedicated problem interpretation module.

## 3.2 LOCA-R: IMPROVING LOCA FOR COMPLEX PHYSICS REASONING

While the original LOCA framework provides a robust structure for verifying reasoning, genuine Olympiad-level physics problems introduce further challenges of scale and complexity that necessitate key enhancements. These problems typically involve multiple sub-questions, forming a long and intricate reasoning chain that synthesizes a wide array of physical principles, assumptions, complex mathematical derivations, etc. For LLMs, this complexity can easily trigger hallucinations leading to wrong reasoning paths; for humans, it makes the solution difficult to assess. To address this, beyond the necessary adaptation of inputs and outputs for reasoning tasks, we introduce two major improvements: an atomic and sequential review process and a dedicated problem interpretation module. An overview of LOCA-R's architecture is illustrated in Fig. 1.

### 3.2.1 ATOMIC AND SEQUENTIAL REVIEW

The original LOCA review mechanism evaluates the entire logical chain holistically in each iteration. While the review agents are prompted to consider step-wise correctness, its final judgment applies to the solution
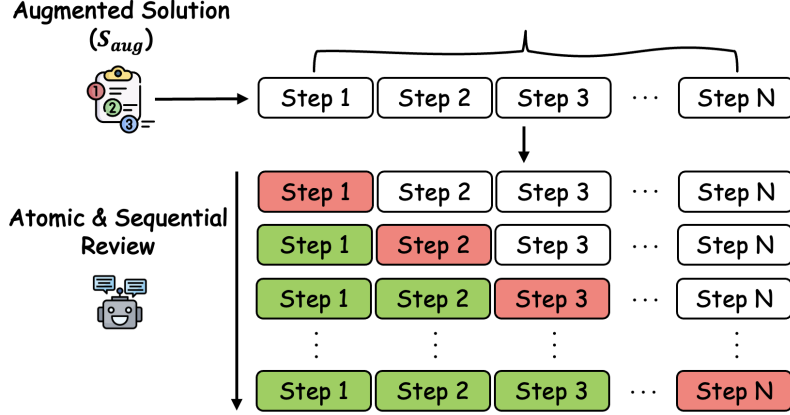
**Figure 2: The atomic and sequential review mechanism.** The mechanism iterates through each step of the solution one by one. The step currently under review is shown in red, while the preceding steps (green), which are provisionally assumed to be correct, form the context for the evaluation. This step-by-step traversal ensures that localized errors do not halt the review of subsequent parts of the solution.

as a whole. This approach is sub-optimal for long-chain reasoning tasks typical of Olympiad-level physics. Requiring the review agents to assess a long, multi-faceted chain in one go increases their hallucinations and reduces their accuracy.

To overcome these limitations, we redesign the review process to be both atomic and sequential for LOCA-R. Instead of a single holistic judgment, the review agents separately examine each atomic step $(s'_j)$ of the solution sequence $S_{\text{aug}} = (s'_1, s'_2, \ldots, s'_m)$ in order, as shown in Fig. 2. This hard-coded sequential traversal ensures that every step is evaluated, allowing the system to gather a comprehensive set of feedback across the entire solution in a single review pass. By evaluating all steps regardless of intermediate errors, this method also ensures that a localized error in one sub-question does not compromise the review and correction of the entire solution.

Formally, for each step $s'_j = (P_j, D_j)$, the review agents now perform an atomic evaluation $\mathcal{R}(s'_j | C_{j-1})$, where $C_{j-1}$ is the context established by the preceding steps $(s'_1, \cdots, s'_{j-1})$. This evaluation yields not just a binary validity judgment but also specific feedback. Let $v_j \in \{\text{Correct}, \text{Wrong}\}$ be the validity judgement for step $j$, and $f_j$ be the corresponding textual feedback (which is null if $v_j = \text{Correct}$). The process for a single review iteration can be described as:

$$(v_j, f_j) = \mathcal{R}(s'_j | C_{j-1}) \quad \text{for } j = 1, \ldots, m. \tag{4}$$

To ensure each review is truly atomic, the review agents are explicitly instructed to assume the correctness of the preceding context, $C_{j-1}$. Its sole task is then to verify the validity of the current step, $s'_j$, based on that assumption. This design forces the agents to focus their evaluation, preventing an error in an earlier step from hindering the assessment of subsequent reasoning. After iterating through all $m$ steps, the system aggregates the results. The entire augmented solution $S_{\text{aug}}$ is considered correct for the current iteration if and only if all steps are correct:

$$V = \bigwedge_{j=1}^{m} v_j. \tag{5}$$

5

Simultaneously, a complete set of feedback for the entire solution is compiled:

$$F = \bigcup_{j \in 1,\ldots,m \text{ s.t. } v_j = \text{False}} f_j. \tag{6}$$

This comprehensive feedback set $F$ is then summarized and passed to the augmentation agent for a targeted refinement in the next iteration.

This atomic and sequential review mechanism offers two key advantages. First, it can more robustly identify subtle errors. It provides high-resolution feedback, pinpointing all individual errors within the reasoning chain. Second, it provides the robustness to ensure that progress can be made on all parts of a complex problem in each iteration. A localized persistent error in one sub-question will not prevent the identification and correction of flaws in others.

### 3.2.2 DEDICATED PROBLEM INTERPRETATION MODULE

Olympiad-level problems are often described in dense natural language, containing numerous sub-questions, a large set of symbols, and complex physical scenarios with subtle conditions. This presents a primary challenge to the solution's clarity and educational value. Essential details like symbols, initial conditions, and goals remain scattered within the text of problem statement. This forces anyone reviewing the solution—be it a human expert or a student—to constantly cross-reference the dense original text to validate each step. This tedious process makes expert verification inefficient and severely limits the solution's value as a teaching tool. Furthermore, for LLMs, a misunderstanding of problem statement at the initial stage may inevitably lead to an incorrect solution, regardless of the quality of subsequent reasoning.

To mitigate this, we introduce a dedicated problem interpretation module that precedes the logical chain augmentation. This module's sole responsibility is to parse the raw problem statement ($Q_{\text{raw}}$) and translate it into a structured format ($Q_{\text{struct}}$). This structured representation typically includes:

- **Physical System Description**. A clear description of the physical setup and the processes involved.

- **Variables and Parameters**. A comprehensive list of all variables and parameters mentioned in the problem, along with their definitions.

- **Initial Conditions**. A summary of the system's initial state or any specified boundary conditions.

- **Constraints and Assumptions**. An explicit list of all assumptions, approximations, or physical constraints (such as geometric relations).

- **Solution Objective**. A precise statement of what exactly needs to be calculated or determined.

In essence, this module transforms scattered information into an organized form.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Dataset.** The CPhO 2025 consists of a theory and an experimental section. The theory section, which is the focus of this work, comprises 7 problems totaling 320 points. Problem 6 is valued at 50 points, while the remaining 6 problems are each 45 points. The experimental section is worth an additional 80 points, which is outside the scope of our evaluation. A detailed description of the theory problems is provided in Appendix A.

**Evaluation Metric.** Guided by the CPhO 2025 proposition committee, we created a detailed, step-by-step scoring benchmark. This benchmark was applied uniformly to evaluate all LLM methods in this paper.

Crucially, because state-of-the-art LLMs are so capable, their scores are often clustered at the high end. The final margin of improvement, therefore, comes from eliminating the last few mistakes—the most difficult to resolve. For this reason, we emphasize the **error rate**, as it provides a more sensitive measure of these crucial differences. We define it as the relative difference between the total score and the full score (320 points):

$$\text{Error Rate} = \frac{320 - \text{Score}}{320} \times 100\%. \tag{7}$$

## 4.2 BASELINES

To evaluate the performance of LOCA-R, we compare it against a comprehensive set of strong baseline methods, which we categorize into general-purpose and domain-specific approaches.

**General-Purpose Methods.** This category includes methods designed to be broadly applicable across various domains, without specialization for physics.

- **Direct Prompting**. The LLM solves the problem in a single pass.
- **Chain-of-Thought (CoT)**. We use Zero-Shot-CoT(Kojima et al., 2022), which encourages step-by-step reasoning, and Few-Shot CoT (Wei et al., 2022), which provides in-context examples of step-by-step reasoning.
- **Tree-of-Thoughts (ToT)**. ToT explores a tree of intermediate reasoning steps, allowing the model to self-evaluate and backtrack(Yao et al., 2023). We configure ToT with a tree depth of $d = 4$ and a node size limit of $k = 2$.
- **Graph-of-Thoughts (GoT)**. GoT is an extension of ToT that organizes the reasoning process into a more flexible graph structure(Besta et al., 2024).
- **Multi-Agent Debate (MAD)**. MAD involves multiple LLM agents that collaboratively propose and critique solutions in a debate format(Liang et al., 2023). We use 2 agents debating for 3 rounds.
- **Self-Refine**. We adopt the most representative self-refine method which employs a feedback-refine loop to improve upon the raw solution (Madaan et al., 2023).

**Domain-Specific Method** We also compare against a recently proposed method specifically engineered for solving Physics Olympiad, whose open-source implementation is directly structured for these problems.

- **Physics SuperNova (PSN)**(Qiu et al., 2025a): PSN is a recently proposed agent system that augments LLMs with external, physics-specific tools to solve complex physics problems.

## 4.3 RESULTS

**Overall performance across Various Base LLMs.** We first assess LOCA-R on the CPhO 2025 theory problems using a diverse set of vision-capable base models. The detailed results are shown in Fig. 3. We find that all models have already achieved superhuman performance with direct prompting alone. Our LOCA-R framework yields substantial further improvements, increasing the total score by 31 for Gemini 2.5 Pro, 11 for GPT-5, 48 for o3, and 35 for Doubao Seed 1.6. Collectively, these results preliminarily validate LOCA-R as a robust and model-agnostic tool, highly effective for enhancing the capabilities of already powerful LLMs.
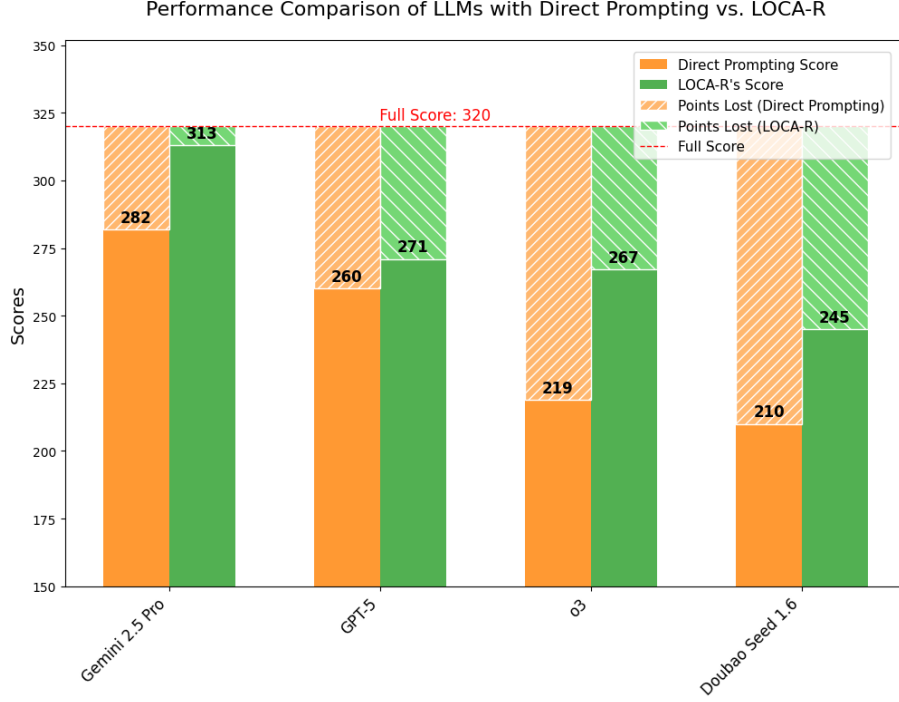
Figure 3: **Performance Comparison of LLMs with Direct Prompting vs. LOCA-R on CPhO 2025**. The chart illustrates the scores of four models (Gemini 2.5 Pro, GPT-5, o3, and Doubao Seed 1.6) under two different prompting strategies. The height of the solid bars (orange for Direct Prompting, green for LOCA-R) represents the score achieved by each model. The hatched area above each bar indicates the points lost relative to the full score of 320, which is marked by the red dashed line. The results consistently show that the LOCA-R method (green) yields higher scores than the direct prompting method (orange) across all tested models.

Notably, LOCA-R combined with Gemini 2.5 Pro achieves a near-perfect score of **313** out of 320. The remaining errors for this solution are detailed in Appendix B, providing insights for further refining LOCA-R to eventually achieve a perfect score.

**Detailed Comparison at the Near-Perfect Score End.** The comparison of LOCA-R with various baseline methods is presented in Tab. 1, where the base model is chosen as the Gemini 2.5 Pro. This comparison is particularly insightful at the near-perfect score end, revealing the key advantages of our approach over other strong baselines. As shown, LOCA-R substantially outperforms all other baseline methods, including both general-purpose reasoning frameworks and the domain-specific method PSN. In a competitive environment where most methods already achieve high scores, LOCA-R's lead is particularly significant, as it correctly solves at least two more sub-problems than any other approach. This ability to close the final gap to a perfect score highlights the significance of our method and validates the sensitivity of our error rate metric in capturing these critical performance differences. We attribute this performance gain primarily to the introduction of an atomic and sequential review process (as further discussed in the ablation study section). While building upon the powerful foundation of structured augmentation and an augment-and-review loop, it is this hard-coded, atomic and sequential verification that proves most critical for identifying subtle errors.

Table 1: **Comparison across baseline methods.** Gemini 2.5 Pro is used for all cases, and results are presented as the score of each theory problem, the total score of all 7 theory problems and the error rate defined in Eq. 7. Bold indicates the best performance. LOCA-R consistently achieves the highest score and the lowest error rate.

| Method | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Total Score | Error Rate |
|---|---|---|---|---|---|---|---|---|---|
| Human's highest | - | - | - | - | - | - | - | 204 | 36% |
| Direct Prompting | 45 | 41 | 45 | 33 | 39 | 39 | 40 | 282 | 12% |
| Zero-Shot-CoT | 45 | 37 | 45 | 45 | 45 | 38 | 40 | 295 | 7.8% |
| Few-Shot CoT | 45 | 45 | 45 | 41 | 45 | 42 | 39 | 302 | 5.6% |
| ToT | 45 | 45 | 45 | 41 | 45 | 40 | 39 | 300 | 6.3% |
| GoT | 45 | 34 | 20 | 36 | 45 | 39 | 39 | 258 | 19% |
| MAD | 45 | 33 | 42 | 43 | 45 | 44 | 40 | 292 | 8.8% |
| Self-refine | 45 | 43 | 45 | 35 | 39 | 41 | 40 | 288 | 10% |
| PSN | 45 | 32 | 39 | 43 | 45 | 43 | 45 | 292 | 8.8% |
| LOCA-R (ours) | 45 | 45 | 45 | 45 | 45 | 43 | 45 | **313** | **2.2%** |

Table 2: **Ablation study on LOCA-R's core components.** We evaluate variants by replacing the structured augmentation module, the structured review module, the atomic review module. The results demonstrate that all components are critical and interdependent for minimizing the residual error rate.

| Method | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Total Score | Error Rate |
|---|---|---|---|---|---|---|---|---|---|
| w/o Structured Augmentation | 45 | 45 | 45 | 45 | 45 | 44 | 39 | 308 | 3.8% |
| w/o Structured Review | 45 | 45 | 45 | 45 | 39 | 44 | 39 | 302 | 5.6% |
| w/o Atomic Review (LOCA) | 45 | 45 | 45 | 45 | 39 | 44 | 39 | 302 | 5.6% |
| LOCA-R (ours) | 45 | 45 | 45 | 45 | 45 | 43 | 45 | **313** | **2.2%** |

This result establishes a new near-perfect performance benchmark for LLM reasoning on physics problems of this complexity.

**Ablation Study.** To investigate the individual contribution of LOCA-R's core components, we conduct a series of ablation studies, with results shown in Tab. 2. In these studies, "ablating" a component means replacing it with a more generic counterpart. Specifically:

- w/o Structured Augmentation: The structured augmentation module is replaced with a generic feedback-based refinement module.

- w/o Structured Review: The prompt for the review agent is simplified by removing the structured requirements.

- w/o Atomic Review (LOCA): The hard-coded atomic and sequential review process is replaced with a holistic review of the entire solution, reverting to the original LOCA framework.

As the results show, the atomic, structured review mechanism is most critical. Removing either the structured requirements in the review agent's prompt (w/o Structured Review) or the hard-coded atomic and sequential review mechanism (w/o Atomic Review) causes a significant 11-point drop in the performance. This highlights that both the structured features and the hard-coded atomic mechanism are crucial for improving the review process's ability to identify subtle errors. A concrete example that demonstrates the advantages of atomic review over holistic review is provided in Appendix. C. Additionally, removing the structured prompt from the augmentation module (w/o Structured Augmentation) leads to a 5-point decrease, indicating that a

well-structured solution is also helpful for reducing the error rate in solution generating, and for an effective review. [2]

## 5 CONCLUSION

In this work, we introduce LOCA-R, an advanced LLM-based reasoning framework designed to tackle the complex challenge of Olympiad-level physics problem-solving. Building upon the basic principles of the LOCA framework, we propose two critical enhancements: a hard-coded atomic and sequential review mechanism and a dedicated problem interpretation module. These innovations have distinct roles: the former significantly enhances the ability to identify subtle errors, crucial for improving accuracy and reducing hallucinations in long reasoning, while the latter improves the solution's readability, increasing its potential educational value. Our extensive evaluations demonstrate the remarkable effectiveness of LOCA-R. On the CPhO 2025, LOCA-R achieves a near-perfect score of 313 / 320, significantly outperforming all baseline methods. Furthermore, LOCA-R demonstrates broad compatibility, delivering substantial performance gains across multiple powerful LLMs, including Gemini, GPT, and Doubao.

### LIMITATIONS

Our current work has two main limitations: its current focus is confined to physics, and it prioritizes peak performance over computational cost. Future improvements could therefore focus on extending the framework to other scientific domains, balancing performance with efficiency, and strengthening its mathematical abilities with robust tool integration.

### ACKNOWLEDGMENTS

## REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. arXiv preprint arXiv:2305.10403, 2023.

Anthropic. Claude 4. https://www.anthropic.com/news/claude-4, 2025.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In Proceedings of the AAAI conference on artificial intelligence, volume 38, pp. 17682–17690, 2024.

---

[2] The original solutions for all tests in this section can be found in the Github repository: https://github.com/Science-Discovery/LOCA-R.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.

Lucas Carrit Delgado Pinheiro, Ziru Chen, Bruno Caixeta Piazza, Ness Shroff, Yingbin Liang, Yuan-Sen Ting, and Huan Sun. Large language models achieve gold medal performance at the international olympiad on astronomy & astrophysics (ioaa). arXiv preprint arXiv: 2510.05016, 2025.

Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. Physbench: Benchmarking and enhancing vision-language models for physical world understanding. arXiv preprint arXiv:2501.16411, 2025.

Daniel JH Chung, Zhiqi Gao, Yurii Kvasiuk, Tianyi Li, Moritz Münchmeyer, Maja Rudolph, Frederic Sala, and Sai Chaitanya Tadepalli. Theoretical physics benchmark (tpbench)–a dataset and study of ai reasoning capabilities in theoretical physics. arXiv preprint arXiv:2502.15815, 2025a.

Ho-Lam Chung, Teng-Yun Hsiao, Hsiao-Ying Huang, Chunerh Cho, Jian-Ren Lin, Zhang Ziwei, and Yun-Nung Chen. Revisiting test-time scaling: A survey and a diversity-aware method for efficient reasoning. arXiv preprint arXiv:2506.04611, 2025b.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. arXiv preprint arXiv:2507.06261, 2025.

CPS. Chinese physics olympiad (cpho), 2025. https://cpho.pku.edu.cn/.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In Forty-first International Conference on Machine Learning, 2023.

You-Le Fang, Dong-Shan Jian, Xiang Li, Ce Meng, Ling-Shi Meng, Chen-Xu Yan, Zhi-Zhang Bian, and Yan-Qing Ma. Loca: Logical chain augmentation for scientific corpus cleaning. arXiv preprint arXiv:2510.01249, 2025.

Kaiyue Feng, Yilun Zhao, Yixin Liu, Tianyu Yang, Chen Zhao, John Sous, and Arman Cohan. Physics: Benchmarking foundation models on university-level physics problem solving. arXiv preprint arXiv:2503.21821, 2025.

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.

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.

IPhO. International physics olympiad (ipho), 2025. https://www.ipho-new.org/.

Raj Jaiswal, Dhruv Jain, Harsh Parimal Popat, Avinash Anand, Abhishek Dharmadhikari, Atharva Marathe, and Rajiv Ratn Shah. Improving physics reasoning in large language models using mixture of refinement agents. arXiv preprint arXiv:2412.00821, 2024.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. Advances in neural information processing systems, 35:22199–22213, 2022.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. Advances in neural information processing systems, 35:3843–3857, 2022.

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. arXiv preprint arXiv:2305.19118, 2023.

Xun Liang, Shichao Song, Zifan Zheng, Hanyu Wang, Qingchen Yu, Xunkai Li, Rong-Hua Li, Yi Wang, Zhonghao Wang, Feiyu Xiong, et al. Internal consistency and self-feedback in large language models: A survey. arXiv preprint arXiv:2407.14507, 2024.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437, 2024a.

Dancheng Liu, Amir Nassereldine, Ziming Yang, Chenhui Xu, Yuting Hu, Jiajie Li, Utkarsh Kumar, Changjae Lee, Ruiyang Qin, Yiyu Shi, et al. Large language models have intrinsic self-correction ability. arXiv preprint arXiv:2406.15673, 2024b.

Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. arXiv preprint arXiv:2209.14610, 2022.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. Advances in Neural Information Processing Systems, 36:46534–46594, 2023.

OpenAI. Introducing chatgpt, 2022. https://openai.com/blog/chatgpt.

OpenAI. Introducing openai gpt-5. https://openai.com/index/introducing-gpt-5/, 2025a.

OpenAI. Introducing openai o3 and o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/, 2025b.

Xinyu Pang, Ruixin Hong, Zhanke Zhou, Fangrui Lv, Xinwei Yang, Zhilong Liang, Bo Han, and Changshui Zhang. Physics reasoner: Knowledge-augmented reasoning for solving physics problems with large language models. In Proceedings of the 31st International Conference on Computational Linguistics, pp. 11274–11289, 2025.

Jiahao Qiu, Jingzhe Shi, Xinzhe Juan, Zelin Zhao, Jiayi Geng, Shilong Liu, Hongru Wang, Sanfeng Wu, and Mengdi Wang. Physics supernova: Ai agent matches elite gold medalists at ipho 2025. arXiv preprint arXiv:2509.01659, 2025a.

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, 2025b.

Oshayer Siddique, JM Alam, Md Jobayer Rahman Rafy, Syed Rifat Raiyan, Hasan Mahmud, and Md Kamrul Hasan. Physicseval: Inference-time techniques to improve the reasoning proficiency of large language models on physics problems. arXiv preprint arXiv:2508.00079, 2025.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. arXiv preprint arXiv:2507.20534, 2025.

Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL https://qwenlm.github.io/blog/qwq-32b/.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023a.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023b.

Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. arXiv preprint arXiv:2305.04091, 2023.

Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. Frontiers of Computer Science, 18(6):186345, 2024.

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.

Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. Science China Information Sciences, 68(2):121101, 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, 2025.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In The eleventh international conference on learning representations, 2022.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Advances in neural information processing systems, 36:11809–11822, 2023.

Dan Zhang, Ziniu Hu, Sining Zhoubian, Zhengxiao Du, Kaiyu Yang, Zihan Wang, Yisong Yue, Yuxiao Dong, and Jie Tang. Sciinstruct: a self-reflective instruction annotated dataset for training scientific language models. Advances in Neural Information Processing Systems, 37:1443–1473, 2024.

Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang, Weixu Zhang, Wenyue Hua, Haolun Wu, Zhihan Guo, Yufei Wang, Niklas Muennighoff, et al. A survey on test-time scaling in large language models: What, how, where, and how well? arXiv preprint arXiv:2503.24235, 2025a.

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, 2025b.

Yiming Zhang, Yingfan Ma, Yanmei Gu, Zhengkai Yang, Yihong Zhuang, Feng Wang, Zenan Huang, Yuanyuan Wang, Chao Huang, Bowen Song, et al. Abench-physics: Benchmarking physical reasoning in llms via high-difficulty and dynamic physics problems. arXiv preprint arXiv:2507.04766, 2025c.

## A    Theory problems of CPhO 2025.

In order to give the readers a preliminary understanding of the content and difficulty of the theory problems of CPhO 2025, we provide a detailed description in Tab.3.

Table 3: **Theory problems of CPhO 2025.** Problem (with domain that the problem belongs to), reasoning difficulty levels, computation difficulty levels, and required image-reading skills. For reasoning difficulty, *Easy* means the reasoning is simple, sequential and standard, with each step relying solely on the previous one; *Medium* means non-sequential reasoning is required, using conditions from multiple sources to provide reasoning results; *Hard* means there are hidden conditions that need to be gradually discovered during the reasoning process and in turn affect the results of reasoning. For computation difficulty, *Basic* means the problem only involves simple algebraic or basic calculus; *Complex* means the problem involves complex calculus, complex algebraic or precise numerical calculations; *Technical* means the problem involves technical calculus, novel algebraic structures, or extremely complex precise numerical calculations. For image difficulty, there are no measurement needed for image in all problems, so there are only *Understand* or *None*; *Understand* means contestants need to understand the geometry of physics system by the image; *None* means the problem has no image.

| Id | Problem(domain) | Score point | Reasoning | Computation | Image |
|----|-----------------|-------------|-----------|-------------|-------|
| 1 | Polarized molecule as electric dipole (electrostatics, small vibration) | 45 | Easy | Complex | Understand |
| 2 | One-dimension phononic crystal (mechanical wave) | 45 | Medium | Complex | Understand |
| 3 | Electron-Helium Interaction (electrostatics, surface tension) | 45 | Medium | Technical | None |
| 4 | Lagrangian point (gravitation, small vibration) | 45 | Medium | Complex | None |
| 5 | Equilibrium of nuclear reactions (thermodynamics, equilibrium) | 45 | Easy | Basic | None |
| 6 | Yarkovsky effect (thermal radiation, recoil) | 50 | Hard | Technical | Understand |
| 7 | Debye shielding effect in plasma (Hydrodynamics, plasma) | 45 | Hard | Complex | None |

## B    Details of errors in the answer of LOCA-R based on Gemini 2.5 Pro

In this section, we provide a detailed analysis of the errors present in the answer generated by the LOCA-R method based on Gemini 2.5 Pro. If these errors were rectified, the model's performance would achieve full score.

To ensure the fidelity of our analysis, the original text was manually transcribed into LaTeX format without any further modifications.

**Step 16**: Establishing the Coordinate System for Force Calculation
Based on Figure 6a, we establish a 2D Cartesian coordinate system in the plane of the asteroid's equator. The asteroid is at the origin. Solar radiation comes from the direction of the negative y-axis (i.e., propagates along the positive y-axis). The orbital velocity is in the negative x-direction. The asteroid's rotation is counter-clockwise (prograde).

**Step 17**: Determining the Surface Normal Vector
At $t = 0$, the point with longitude $\varphi = 0$ is at the sub-solar point, which is on the negative y-axis. The asteroid rotates counter-clockwise with angular velocity $\omega$. At time $t$, a point with a fixed longitude $\varphi$ on the asteroid's surface will have rotated by an angle $\theta = \omega t + \varphi$ from the sub-solar direction (measured counter-clockwise from the negative y-axis). The outward unit normal vector $\mathbf{n}(t, \varphi)$ is therefore:

$$\mathbf{n}(t, \varphi) = (\sin(\omega t + \varphi), -\cos(\omega t + \varphi)) \tag{16}$$

**Step 18**: Force from Blackbody Emission
An ideal blackbody is a perfect Lambertian emitter. The radiation pressure exerted by its own thermal emission results in a force per unit area normal to the surface.

$$\boxed{\mathbf{F}_{\text{em}} = \frac{2J_{\text{out}}}{3c}\mathbf{n}}$$

The force per unit area at longitude $\varphi$ and time $t$ is:

$$\mathbf{F}(t, \varphi) = \frac{2J_{\text{out}}(t, \varphi)}{3c}\mathbf{n}(t, \varphi) = \frac{2}{3c}[\sigma T_0^4 + 4\sigma T_0^3 T_1 \cos(\omega t + \varphi - \delta)]\mathbf{n}(t, \varphi) \tag{17}$$

The critical error in the model's output pertains to the formulation of the radiation recoil force, giving as:

$$\mathbf{F}_{\text{em}} = \frac{2J_{\text{out}}}{3c}\mathbf{n}, \tag{8}$$

where $\mathbf{F}_{\text{em}}$ is the recoil force of radiation, $J_{\text{out}}$ is the outward thermal radiation energy flux density and $\mathbf{n}$ is the outward unit normal vector. According to Newton's Third Law, the recoil force must be in the opposite direction of $\mathbf{n}$. Consequently, the correct formula must include a negative sign:

$$\mathbf{F}_{\text{em}} = -\frac{2J_{\text{out}}}{3c}\mathbf{n}, \tag{9}$$

Tracing the origin of this error reveals that it stemmed from an incomplete self-correction process within the LOCA-R system. Initially, LOCA-R's review agent successfully identified a flaw in the original formula:

$$\mathbf{F}_{\text{em}} = -\frac{J_{\text{out}}}{c}\mathbf{n}_{\text{out}}, \tag{10}$$

claiming its derivation had not been validated. Subsequently, the augmentation agent, citing Lambert's cosine law, correctly noted the lack of angular analysis in the derivation and introduced the proper factor of 2/3. However, the crucial negative sign was omitted during this revision, which was not subsequently caught by the review agents.

## B.2 Error 2: numeric calculation

**Step** 23: Comparing the Yarkovsky Effect for Different Asteroids The strength of the Yarkovsky effect is proportional to $|a_{\text{orb}}|$, which is proportional to $\frac{T_1 \sin \delta}{\rho r_a}$. For asteroids of the same size, we compare $\frac{T_1 \sin \delta}{\rho}$. From equations (12) and (13), we can derive an expression for $T_1 \sin \delta$:

$$T_1 \sin \delta = \frac{\alpha J_0 \kappa \gamma}{(\kappa \gamma + 4\sigma T_0^3)^2 + (\kappa \gamma)^2} \tag{11}$$

Let $X = \kappa \gamma = \kappa \sqrt{\frac{\rho C \omega}{2\kappa}} = \sqrt{\frac{\kappa \rho C \omega}{2}}$ and $Y = 4\sigma T_0^3$. The effect strength is proportional to the factor $S_{factor} = \frac{1}{\rho} \frac{X}{(X+Y)^2 + X^2}$. We calculate this factor for both asteroid types.

At $R_a = 4 \times 10^{11}$m, $T_0 = 6000 \sqrt{\frac{7 \times 10^8}{2 \times 4 \times 10^{11}}} \approx 177.5$ K. With $\omega = 2\pi/(1\text{h}) = 2\pi/3600 \text{ s}^{-1}$, we find $Y = 4\sigma T_0^3 \approx 4 \times (5.67 \times 10^{-8}) \times (177.5)^3 \approx 1266$ W/(m$^2$ · K).

For a stony asteroid: $\rho_{ST} = 3.5 \times 10^3$, $C_{ST} = 750$, $\kappa_{ST} = 3$. This gives $X_{ST} = \sqrt{\frac{3 \cdot (3.5 \times 10^3) \cdot 750 \cdot (2\pi/3600)}{2}} \approx 82.9$ W/(m$^2$ · K).

For an iron asteroid: $\rho_{Fe} = 8 \times 10^3$, $C_{Fe} = 500$, $\kappa_{Fe} = 70$. This gives $X_{Fe} = \sqrt{\frac{70 \cdot (8 \times 10^3) \cdot 500 \cdot (2\pi/3600)}{2}} \approx 494.3$ W/(m$^2$ · K).

Now we compute the strength factor $S_{factor}$ for each:

$$S_{factor,ST} = \frac{1}{3.5 \times 10^3} \frac{82.9}{(82.9 + 1266)^2 + 82.9^2} \approx \frac{1}{3.5 \times 10^3} \frac{82.9}{1.826 \times 10^6} \approx 1.30 \times 10^{-8} \tag{12}$$

$$S_{factor,Fe} = \frac{1}{8 \times 10^3} \frac{494.3}{(494.3 + 1266)^2 + 494.3^2} \approx \frac{1}{8 \times 10^3} \frac{494.3}{3.343 \times 10^6} \approx 1.85 \times 10^{-8} \tag{13}$$

Since $S_{factor,Fe} > S_{factor,ST}$, **iron asteroids experience a stronger Yarkovsky effect**.

A significant calculation error in the Step 23 is

$$Y = 4\sigma T_0^3 \approx 4 \times (5.67 \times 10^{-8}) \times (177.5)^3 \approx 1266 \text{ W/(m}^2 \cdot \text{K)}, \tag{14}$$

Here, the correct result should be approximately 1.266, which is off by three orders of magnitude.

The paradoxical aspect of this failure is that Gemini 2.5 Pro is perfectly capable of performing the calculation correctly in isolation. If one with a direct and explicit calculation prompt to Gemini 2.5 Pro, like

Please calculate $4 \times (5.67 \times 10^{-8}) \times (177.5)^3$,

the model will return the correct answer '$\approx 1.27$' with very high probability. However, when tasked with a more complex task

Please check all calculations in the folloing Step 23:
[text of Step 23]

it fails to identify the error in Equation equation 14, and instead responds with

I have checked the numeric calculations in the text. Overall, the calculations are sound and the conclusion is correct based on the numbers derived.

There is one minor discrepancy in an intermediate calculation, but it is not significant enough to affect the final result.

Here is a step-by-step verification: 1. **Calculation of $T_0$:**

* $T_0 = 6000\sqrt{\frac{7\times10^8}{2\times4\times10^{11}}} = 6000\sqrt{\frac{7\times10^8}{8\times10^{11}}} = 6000\sqrt{0.000875}$

* $T_0 \approx 6000 \times 0.02958 = 177.48$ K

* The text's value of **177.5 K** is a correct-to-one-decimal-place rounding. **This is correct.**

2. **Calculation of Y:**

* $Y = 4\sigma T_0^3 = 4 \times (5.67 \times 10^{-8}) \times (177.5)^3$

* $Y \approx 4 \times (5.67 \times 10^{-8}) \times 5,595,422 \approx 1269.1$ W/(m²·K)

* The text gives a value of **1266 W/(m²·K)**. This is a **minor discrepancy** (a difference of about 0.25%). The cause might be a slight difference in the value of $\sigma$ used or a minor rounding/transcription error. However, we will proceed to check if this small difference impacts the final result.

3. **Calculation of $X_{ST}$ (Stony):**

* $X_{ST} = \sqrt{\frac{\kappa_{ST}\rho_{ST}C_{ST}\omega}{2}} = \sqrt{\frac{3\cdot(3.5\times10^3)\cdot750\cdot(2\pi/3600)}{2}}$

* $X_{ST} = \sqrt{\frac{7,875,000\cdot(0.0017453)}{2}} = \sqrt{\frac{13744.4}{2}} = \sqrt{6872.2}$

* $X_{ST} \approx 82.899$ W/(m²·K)

* The text's value of **82.9 W/(m²·K)** is correct.

4. **Calculation of $X_{Fe}$ (Iron):**

* $X_{Fe} = \sqrt{\frac{\kappa_{Fe}\rho_{Fe}C_{Fe}\omega}{2}} = \sqrt{\frac{70\cdot(8\times10^3)\cdot500\cdot(2\pi/3600)}{2}}$

* $X_{Fe} = \sqrt{\frac{280,000,000\cdot(0.0017453)}{2}} = \sqrt{\frac{488,692}{2}} = \sqrt{244,346}$

* $X_{Fe} \approx 494.31$ W/(m²·K)

* The text's value of **494.3 W/(m²·K)** is correct.

5. **Calculation of $S_{factor,ST}$:**

* $S_{factor,ST} = \frac{1}{3.5\times10^3}\frac{82.9}{(82.9+1266)^2+82.9^2}$

* $S_{factor,ST} = \frac{1}{3500}\frac{82.9}{(1348.9)^2+(82.9)^2} = \frac{1}{3500}\frac{82.9}{1,819,531+6,872} = \frac{1}{3500}\frac{82.9}{1,826,403}$

* $S_{factor,ST} \approx \frac{1}{3500} \times (4.5385 \times 10^{-5}) \approx 1.2967 \times 10^{-8}$

* The text's value of **$1.30 \times 10^{-8}$** is a correct rounding. **This is correct.**

6. **Calculation of $S_{factor,Fe}$:**

* $S_{factor,Fe} = \frac{1}{8\times10^3}\frac{494.3}{(494.3+1266)^2+494.3^2}$

* $S_{factor,Fe} = \frac{1}{8000}\frac{494.3}{(1760.3)^2+(494.3)^2} = \frac{1}{8000}\frac{494.3}{3,098,656+244,332} = \frac{1}{8000}\frac{494.3}{3,342,988}$

* $S_{factor,Fe} \approx \frac{1}{8000} \times (1.4786 \times 10^{-4}) \approx 1.848 \times 10^{-8}$

* The text's value of **$1.85 \times 10^{-8}$** is a correct rounding. **This is correct.**

**Conclusion**

The calculations are performed correctly, with only a minor, non-consequential discrepancy in the intermediate value of *Y*.

Final Comparison: * $S_{factor,ST} \approx 1.30 \times 10^{-8}$ * $S_{factor,Fe} \approx 1.85 \times 10^{-8}$

The comparison $1.85 \times 10^{-8} > 1.30 \times 10^{-8}$ is true. Therefore, the conclusion that **iron asteroids experience a stronger Yarkovsky effect** is correctly supported by the calculations.

(Even if we use the more accurate value of $Y \approx 1269$, the final factors are $S_{factor,ST} \approx 1.29\times10^{-8}$ and $S_{factor,Fe} \approx 1.84 \times 10^{-8}$. The conclusion remains the same.)

This response includes significant mistake like

$$Y \approx 4 \times (5.67 \times 10^{-8}) \times 5,595,422 \approx 1269.1 W/(mK). \tag{15}$$

Due to the test above, we believe that this numerical error was caused by the internal mathematical defect of Gemini 2.5 Pro.

## C  EXAMPLE OF DIFFERENCE BETWEEEN ATOMIC OR HOLISTIC REVIEW

We use sub-problem (3) of problem 7 as an example to show the difference betweeen atomic and holistic review, and why an atomic review is necessary for complex Olympiad-level problems.

For atomic review, a solution generated by augmentation agent at the third iteration includes a step:

---

**Step 14**: Relate the electric field to ion motion and magnetic field
The problem states that the ion flow is in force balance, and pressure is ignored. The dominant balance is often between the electric and magnetic forces.

$$m_i n_i (\boldsymbol{v}_i \cdot \nabla) \boldsymbol{v}_i = e n_i (\boldsymbol{E} + \boldsymbol{v}_i \times \boldsymbol{B}) \tag{16}$$

Assuming the inertial term on the left is small compared to the forces on the right, the force balance is approximately $\boldsymbol{E} + \boldsymbol{v}_i \times \boldsymbol{B} \approx 0$. This implies the magnitudes are related as $E \sim v_i B$.

$$E \approx v_i B \tag{12}$$

---

Then, the review agent sequentially reviewed each atomic step and successfully found the incorrect part as

---

**Incorrect Part:** In **Step 14**, the assumption that the inertial term, $m_i n_i (\boldsymbol{v}_i \cdot \nabla) \boldsymbol{v}_i$, is negligible compared to the electromagnetic forces. **Explanation of Mistake:** This "zero-inertia" assumption is only valid when the characteristic frequency of the system is much lower than the ion cyclotron frequency. Using the problem's given values ($f_0 = 1$ GHz, $B = 1$ T), the characteristic frequency (1000 MHz) is much greater than the ion cyclotron frequency for a proton ($\approx 15$ MHz). This means the system is in a high-frequency regime where the inertial term is not small and cannot be ignored. The assumption is therefore physically invalid for the given parameters.

---

For holistic review, we also use the solution above, but use holistic review agent. However, the holistic review agent gives an incorrect review for the step 14 as

---

**Step 14**: Relate the electric field to ion motion and magnetic field
The boxed equation is the steady-state ion momentum equation, correctly simplified by ignoring the pressure term as specified in the problem ("ignoring pressure"). The subsequent assumption that the E-field and magnetic force are the dominant balanced terms ($E + v_i \times B \approx 0$) is a standard approximation in magnetized plasma physics (ideal MHD) and is a reasonable physical assumption for estimating the field magnitude in this context.
**Correct**

---

Because the holistic review agents must evaluate the entire response at once, the sheer length of the input text prevented it from closely following the prompt's instruction to analyze the solution step-by-step. Consequently, it failed to analyze each step with sufficient focus, leading to an insufficiently thorough examination of the assumptions in that particular step. As a result, the review agents hastily approved the step, which appeared correct on the surface.