

The Stochastic Parrot on LLM’s Shoulder: A Summative Assessment of Physical Concept Understanding

Anonymous ACL submission

Abstract

In a systematic way, we investigate a widely asked question: *Do LLMs really understand what they say?*, which relates to the more familiar term *Stochastic Parrot*. To this end, we propose a summative assessment over a carefully designed physical concept understanding task, PHYSICO. Our task alleviates the memorization issue via the usage of grid-format inputs that abstractly describe physical phenomena. The grids represents varying levels of understanding, from the core phenomenon, application examples to analogies to other abstract patterns in the grid world. A comprehensive study on our task demonstrates that: (1) state-of-the-art LLMs lag behind humans by $\sim 40\%$; (2) the stochastic parrot phenomenon is present in LLMs, as they fail on our grid task but can describe and recognize the same concepts well in natural language; (3) our task challenges the LLMs due to intrinsic difficulties rather than the unfamiliar grid format, as in-context learning and fine-tuning on same formatted data added little to their performance. Our data and code will be released for public research.

1 Introduction

Recent years have witnessed remarkable advancements in large language models (LLMs) (Brown et al., 2020; Achiam et al., 2023; Team et al., 2023). Thanks to the substantial model capacity and massive training data, LLMs have achieved new state-of-the-arts on a variety of NLP tasks are even surpassing humans on some of them (Min et al., 2023; Chang et al., 2024). Nowadays the application of LLMs has become widespread, facilitating daily work and life, and profoundly influencing people’s work and lifestyles (Bommasani et al., 2021; Peng et al., 2024; Demszky et al., 2023).

On the other hand, despite the great success of LLMs, many researchers argue that *LLMs may not really understand what they claim they do* (Bender and Koller, 2020; Bender et al., 2021; Bom-

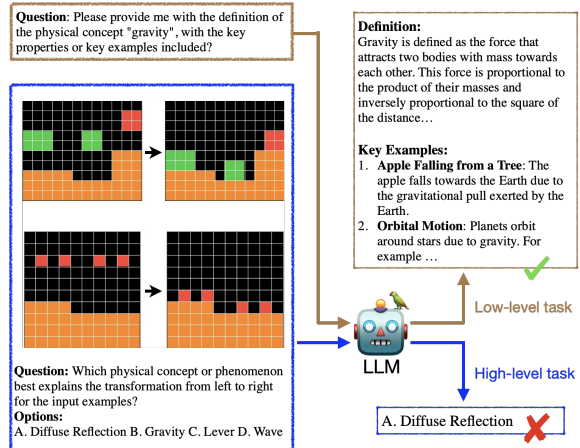


Figure 1: Illustration of a “Stochastic Parrot” by our PHYSICO task consisting of both **low-level** and **high-level** subtasks in parallel. For a concept *Gravity*, an LLM can generate its accurate description in natural language, but cannot interpret its grid-format illustration.

masani et al., 2021; Mitchell and Krakauer, 2023) due to their strong memorization ability. In particular, Bender et al. (2021) questioned whether LLMs are just *Stochastic Parrots* that repeat words based on correlations without true understanding. This argument has been acknowledged by many research papers and dozens of them even include this term in their titles.¹ Unfortunately, to our best knowledge, there are no quantitative experiments to verify the stochastic parrot phenomenon in LLMs. Existing studies indicate that LLMs may fail on one particular challenging task (Chakrabarty et al., 2022; Shapira et al., 2023; Hessel et al., 2023; Tong et al., 2024), but they do not demonstrate that LLMs claimed to understand those tasks by providing a controlled and paired evidence.

This paper aims to provide quantitative evidence to validate the argument of stochastic parrot in LLMs. To this end, from the perspective of educational and cognitive psychology, we first employ the approach of summative assessment (Black

¹https://scholar.google.com/scholar?hl=en&as_sdt=0%2C5&q=llms+are+stochastic+parrot&btnG=

and Wiliam, 1998a,b) to measure understanding in LLMs. Its key idea is to design various tasks that test different understanding levels regarding a specific concept. Following the principle of Bloom’s taxonomy (Armstrong, 2010; Krathwohl, 2002), we design tasks that reflect different levels of understanding. Consequently, we develop PHYSICO, a task designed to assess understanding of basic physical concepts from high school such as *Gravity*. Our focus on physical concepts stems from both their fundamental relevance to important topics of world models and embodied systems (Savva et al., 2019; Duan et al., 2022; Xiang et al., 2023), and their rich denotations and connotations that enable effective design of summative assessment tasks.

Specifically, PHYSICO includes two subtasks corresponding to two coarse levels of understanding in Bloom’s taxonomy, as shown in Figure 1. One is the low-level understanding subtask in the natural language format, aimed at measuring the remembering (or memorization) ability of LLMs. The other involves the same concepts but in an abstract representation format inspired by (Chollet, 2019), which is designed to measure the high-level understanding beyond remembering of LLMs.

We conduct comprehensive experiments on PHYSICO with representative open-source and commercial LLMs. We obtain two key findings: 1) LLMs perform perfectly on the low-level understanding subtask (>95% in Accuracy) but lags behind humans by a large margin (~40% in Accuracy) on the high-level subtask, which verifies the stochastic parrot phenomenon in LLMs. 2) Further analysis shows that our high-level subtask challenges LLMs due to the intrinsic difficulty of deep understanding rather than the unfamiliar format.

This paper makes the following contributions:

- We introduce a psychology-appealing approach (summative assessment) to measure the understanding of LLMs.
- To fulfil summative assessment, we propose a challenging task PHYSICO to evaluate LLMs.
- Based on PHYSICO, we provide a quantitative experiment to successfully verify the stochastic parrot phenomenon in LLMs.

2 Measuring Concept Understanding via Summative Assessment

It is intrinsically challenging to measure the extent to which LLMs understand a sentence or concept. Indeed, Bender and Koller (2020) provide a defini-

tion of "understanding" from a linguistic perspective, but this definition depends on another abstract and unmeasurable term, "meaning". Therefore, even with this definition, accurately measuring "understanding" remains elusive.

In general, "understanding" is not only a term in linguistics but also in educational and cognitive psychology. Hence, we approach the measurement of whether LLMs understand a concept from an educational and cognitive perspective, using **summative assessment** (Black and Wiliam, 1998a). The key idea is illustrated by the following example.

A Motivating Example Suppose a middle school physics teacher is explaining the concept of "Gravity" to students. How can the teacher know whether a student truly understands this concept? In practice, the teacher would design a series of questions specifically related to gravity to assess comprehension, e.g., the properties like inverse square law and examples like orbital motions. If a student struggles to answer many of these questions, the teacher may conclude that the student does not understand the concept well or has a poor grasp of it.

Summative Assessment Summative assessment, which is widely used by educators, is an appealing strategy to evaluate students’ understanding and knowledge acquisition in educational and cognitive psychology (Black and Wiliam, 1998a,b; Harlen and James, 1997). In this paper, we extend it from evaluating humans to evaluating machines.

Assume \mathcal{S} denotes an intelligent system and \mathcal{C} is a specific concept. To evaluate the extent how \mathcal{S} understands the concept \mathcal{C} , our summative assessment includes the following two steps:

- *Task design towards \mathcal{C}* : design several concept understanding tasks, each of which consists of a lot of questions manually created towards understanding the concept \mathcal{C} .
- *Evaluating \mathcal{S}* : ask \mathcal{S} to answer the questions from the tasks and calculate its accuracy.

Requirements for Validity The success (validity) of the proposed evaluation approach highly depends on the task design (Black and Wiliam, 1998a,b). For example, if the questions are too easy, even a weak system could answer them correctly. This leads to an overestimation of the system’s understanding capabilities, making the assessment ineffective. To ensure good validity, we adhere to the principles outlined in summative assessment (Black and Wiliam, 1998a,b) for task design:

- *Alignment with evaluating objectives*: the questions should be related to the targeted concept, and should measure the specific knowledge about the targeted concept.
- *Different difficulty levels*: the questions should be with different difficulty levels from easy to difficult level, to ensure that the evaluation results have distinctiveness for different systems.
- *Variety*: the questions should reflect various understanding aspects of the targeted concept; addressing both its denotation and connotation.

3 Task Design and Dataset Construction

3.1 Task Design Principle

We borrow the idea of Bloom’s taxonomy (Krathwohl, 2002; Armstrong, 2010) from education research to fulfill the requirements for task design in Section 2, so as to ensure the assessment validity. Bloom’s taxonomy offers an ideal principle to these requirements with an ordering of six cognitive skills (from low to high level) for knowledge understanding: *Remembering, Understanding, Applying, Analyzing, Evaluating and Creating*.

Generally, it is nontrivial to strictly follow this principle since there is no clear boundary among the last four skills of understanding. As a result, we group the last four high-level skills into one and consider the following two levels of understanding:

- *Low-level Understanding*: covering the two lowest-level skills in Bloom’s taxonomy, *i.e.*, retrieving relevant knowledge from long-term memory and rephrasing in one’s own words.
- *High-level Understanding*: covering the aspects for understanding the knowledge beyond memorization, *e.g.*, applying the knowledge to explain a physical phenomenon, analyzing a concrete property of a concept in a generalized and abstract manner,² and explaining phenomena by connecting different concepts.

Based on the two levels of understanding, we design PHYSICO task for summative assessment.

3.2 Our PHYSICO Task

PHYSICO is essentially a physical concept understanding task, which primarily targets on 42 physical concepts or phenomena: *e.g.*, *gravity, light reflection, acceleration, buoyancy, inertia, etc* (see Appendix A for the full list). Our focus on physical concepts is motivated by two main reasons: 1)

²For example, the flow of electric current can be abstracted as *moving* from high potential to low potential.

understanding physical concepts is critical for intelligent systems to interact with the world, which is ultimate goal of embodied AI (Savva et al., 2019; Duan et al., 2022; Xiang et al., 2023); 2) designing tasks centered around physical concepts allows us to more easily control different levels of understanding and ensure the diversity of each concept.

For each physical concept, PHYSICO involves both low-level understanding subtasks and high-level subtasks, following our task design principles.

3.2.1 Low-level Understanding Subtasks

Physical Concept Selection Subtask To evaluate whether an LLM possesses the knowledge of our included concepts, we design a task for LLMs to recognize a concept from its corresponding Wikipedia definition. We manually masked the synonyms of the concept with placeholder [PHENOMENON]. Additionally, highly relevant entities were masked as [MASK] to alleviate short-cuts. For example, in the definition of *Gravity*, the terms “gravity” and “gravitation” were masked as [PHENOMENON], while “Isaac Newton” was masked as [MASK]. Detailed process is described in Appendix B. We then present the models with four choices for concept identification, consistent with the following high-level subtasks.

Physical Concept Generation Subtask As the second subtask, we directly ask the LLMs to generate the description of a concept with its core properties and representative examples. For instance, the concept *Gravity* was described as “*a force that pulls objects with mass towards each other*”, followed by the example “*an apple falls to the ground*” as shown in Figure 1. We then evaluate the quality of the description and its coverage of knowledge required by our PHYSICO. This provides a quantitative measure of the knowledge LLMs can recall and rephrase in the context of our assessment.

3.2.2 High-level Understanding Subtasks

The low-level subtasks are depicted in natural language thus are likely to be remembered by the LLMs due to their extensive training data. To assess whether the LLMs possess a deep understanding of the knowledge, we require the subtasks that can 1) represent the high-level understanding skills; 2) avoid the effects of memorization.

The Abstraction and Reasoning Corpus (ARC) (Chollet, 2019) provides a compelling way by using grids (or matrices) instead of texts to represent a concept. While the LLMs have seen

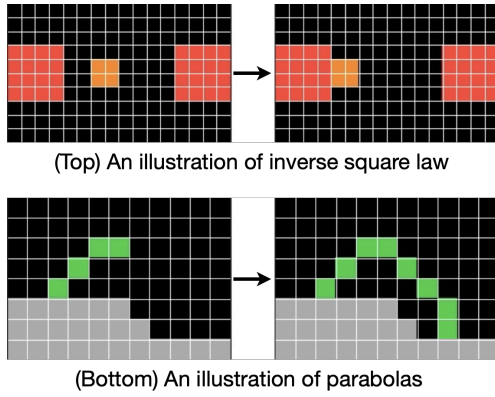


Figure 2: Examples of input-output grid representation labeled as *Gravity*, with increasing difficulty levels.

matrices during pre-training, the data is less likely to be correlated to physical concepts. We hence adopt this idea to represent our subtask as abstract representations in the grid world that associate to the key properties of a physical concept.

The PHYSICO-CORE Set Our first subtask aims to cover the core properties or most representative examples/applications of the assessed concepts. To ensure our set remains generally comprehensible to humans, we maintain a high school-level difficulty and selected 27 common physical concepts within the curriculum. To enhance the diversity and richness, five annotators have labeled multiple core aspects of each concept. For example, the annotated core aspects of *Gravity* include *attraction between two bodies*, *motion on an inclined plane*, *objects falling to grounds* and *orbital motions*.

For each aspect of a concept, the annotator is asked to draw several pairs of abstract grid representations. The aspect of the concept is guaranteed to be illustrated by the pair, such that it explains the transformation from the input to the output. For example, Figure 1 forms a direct abstract visualization of the *Gravity* concept from textbooks, *i.e.*, *apple falling from a tree*. This results in 600 paired instances for the 27 concepts.

Figure 2 presents two examples from this subtask that delve deeper into the concept of *Gravity* compared to Figure 1. The top example demonstrates an application of the *inverse square law of gravity*. The bottom one presents a parabola, linking the knowledge of *gravity* to *inertia*. These examples demonstrate the difficulty of inferring their ground-truth labels solely by recalling the concept of *Gravity* without high-level understanding skills.

The PHYSICO-ASSOCIATIVE Set Many instances in the original ARC dataset can be solved

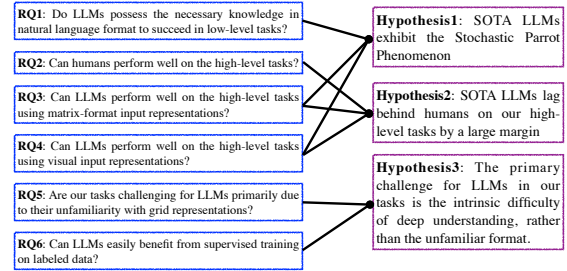


Figure 3: Overview of the research questions answered in our study and their relationships.

via association or analogy to physical concepts. Therefore, as a second source of subtasks, we ask annotators to manually pick input-output grids from ARC that can evoke their associations to specific physical concepts and assign these concepts as ground-truth labels. Different from PHYSICO-CORE, we adopt an open-coding schema and allow the inclusion of new concepts during annotation. The annotators have reviewed 500 ARC instances to filter out the required ones. After cross-validation to ensure agreement, it results in a collection of 200 instances with physical concept labels.

This relabelling approach covers additional 15 physical concepts. The resulted subtasks have each example grid represent an abstract aspect of a concept with possible distracting information. Consequently, the resulted task is more subjective hence more challenging than the PHYSICO-CORE Set.

Creation of Classification Tasks We create *four-choice* tasks on the annotated data. Each instance consists of at least 3 unique grid pairs as input examples. This results in 200 instances for PHYSICO-CORE and 200 instances for ASSOCIATIVE respectively. For each instance, we select three additional labels from our concept pool, along with the ground-truth label, as candidate options. We manually avoid ambiguity during the negative sampling. For example, if *Gravity* is the ground-truth, concepts like *Magnet* will not be sampled.

4 Overview of Our Studies

In the following sections, we conduct a series of studies on our PHYSICO tasks. Our studies are organized into six *Research Questions (RQs)*, through which we aim to answer three *Hypotheses* as shown in Figure 3. In summary, we propose to:

(1) Examine the quantitative disparity in LLMs' performances between low-level (RQ 1) and high-level subtasks (RQ 3, RQ 4). This aims to highlight **the existence of stochastic parrot phenomenon** in LLMs' understanding of physical concepts.

(2) Assess the performance gap between LLMs (RQ 3, RQ 4) and humans on our high-level subtasks (RQ 2). This aims to demonstrate that LLMs **fall significantly short of human understanding**.

(3) Investigate the shortcomings of in-context learning and supervised fine-tuning in improving LLMs on our high-level subtasks (RQ 5, RQ 6). This aims to underscore the **intrinsic limitations** of SOTA LLMs in achieving deep understanding.

Experimented Models We use commercial LLMs, including GPT-3.5 (gpt-3.5-turbo-1106), GPT-4 and GPT-4v (gpt-4-turbo-2024-04-09) and GPT-4o (gpt-4o-2024-05-13); and open-source LLMs, including Llama-3 (Llama-3-8B-Instruct) (MetaAI, 2024) and Mistral (mistral-7B-Instruct-v0.2) (Jiang et al., 2023), InternVL-Chat-V1-5 (Chen et al., 2023, 2024)³ and LLaVA-NeXT-34B (Liu et al., 2023a,b). We use the default inference configurations of the LLMs. Considering the randomness, we run each experiment 3 times and compute the average and standard derivation.

5 Validation on Low-Level Subtasks

To illustrate the stochastic parrot phenomenon with PHYSICO, a necessary condition is to ensure the LLMs can perform well on the low-level understanding subtasks, *i.e.*, whether LLMs exhibit strong skills of *recalling* and *describing* the definitions, core properties and representative examples of the physical concepts in our tasks. That is:

RQ 1: Can LLMs perform well on low-level subtasks, *i.e.*, understanding the definitions of physical concepts in natural language?

To answer RQ 1, we evaluate the LLMs’ abilities to comprehend the definitions of these concepts and generate their descriptions and examples in natural language, as defined in Section 3.2.1.

5.1 Concept Selection Subtask

Settings We provide the standard definition of a concept based on Wikipedia with its synonyms masked; then ask the LLMs to identify the concept, under the same four-choice setting throughout the experiments. We evaluate the representative text-only LLMs and compute the accuracy.

Results Table 1 shows that the GPT models perform near perfect on recognition of our physical concepts from standard written definitions. A small

Mistral	Llama-3	GPT-3.5	GPT-4
81.0 \pm 1.3	88.5 \pm 0.7	97.3 \pm 0.3	95.0 \pm 0.9

Table 1: Accuracy on the concept selection subtask.

	Mistral	Llama-3	GPT-3.5	GPT-4
Human	92.6	100	100	100
SP	89.2 \pm 1.6	91.9 \pm 0.6	96.0 \pm 0.4	99.8 \pm 0.2

Table 2: Evaluations on the concept generation subtask, with metrics of Self-Play success and Human evaluation.

number of errors stem from confusing *reflection* or *light imaging* with *refraction*. Open-source models make mistakes on the same concepts but more frequently, showing a misunderstanding of the three concepts. Additionally, the models occasionally fail to follow instruction, and predict a synonym instead of selecting the correct answer.

5.2 Concept Generation Subtask

Settings This subtask evaluates the descriptions LLMs generate for a concept. It is a text generation task, the evaluation of which is in general difficult. Moreover, in our scenario each concept have many different ground-truth examples in its description, thus existing automatic metrics such as BLEU (Papineni et al., 2002) and METEOR (Banerjee and Lavie, 2005) are not capable of accurately measuring the quality. Therefore, we propose an alternative automatic metric via a self-play game as well as human evaluation for this subtask.

- **Automatic evaluation metric via self-play game:** For each generated description of a concept, we mask the synonyms of the concept in it as in the previous selection subtask, and ask the same LLM to identify the concept being described from four options. This metric evaluates the quality of LLMs’ generated concept descriptions in an objective manner.
- **Human evaluation metric:** We ask the annotators to evaluate the quality of the generated descriptions. The evaluation uses binary scores: each description receives a score of 0 if it consists of any factual error on the concept itself or any unfaithful examples,⁴ and a score of 1 otherwise.

Results The results of automatic and human evaluations are shown in Table 2. According to human evaluation, there are no factual errors in the generated descriptions except for Mistral, confirming that our selected concepts rely on basic and widely accepted knowledge. Thought accurate, the

³Best open-source vision-language model at OpenVLM Leaderboard.

⁴For example, if the LLMs generated a wrong year in the description, it is not counted as incorrect physical knowledge.

open-source LLMs sometimes include correct but uncommon facts, *e.g.*, listing single-slit diffraction as an example of *Wave Interference*.

In the self-play test, all LLMs can accurately recognize the physical concepts from the descriptions they wrote by themselves. Combined with the conclusion from human evaluation, it shows the LLMs can generate correct and sufficient information.

Remark We also ask the annotators of our PHYSICO-CORE to evaluate whether the core properties they annotated are covered by the LLMs’ generated descriptions. This corresponds to measuring the recall of the generated descriptions on core properties/examples of concepts from PHYSICO-CORE. The recall rates for GPT-3.5 and GPT-4 are 85.0 and 90.0, respectively. Of course, there are some exceptional examples from PHYSICO-CORE missed in the descriptions. One example is that the LLMs fails to draw the connection between *movable pulley* and the *Lever* concept. Moreover, by manually checking these missed properties and examples, we found that most of them can be recalled if we query the LLMs in a second turn by prompting “Any more core properties or examples?”. This confirms that the LLMs are *aware of* and are *able to recall* the core properties of concepts covered by the PHYSICO-CORE, though some of them may not have the top conditional probabilities of generation.

Conclusion LLMs understand the concepts covered by PHYSICO in natural language format. Notably, we find that the properties and examples annotated in PHYSICO-CORE are *within the LLMs’ knowledge* and are *highly likely to pop up* when the corresponding physical concepts are queried.

6 Experiments on High-Level Subtasks

This section answers the research questions regarding our high-level understanding subtasks.

RQ 2: Can Humans understand the abstract representations?

First of all, we investigate the performance of humans who possess the knowledge required by our PHYSICO. For each instance in our PHYSICO, we asked three independent annotators who were *not* involved in our task design to perform the same classification task presented to the LLMs. The results indicate that our tasks are largely solvable to people with a college-level education. Specifically, on the PHYSICO-CORE tasks, humans achieved

	Models	CORE	ASSOCIATIVE
	Random	25.0	25.0
text-only	GPT-3.5	26.5 \pm 2.5	30.0 \pm 2.5
	GPT-4	41.3 \pm 1.3	38.3 \pm 1.2
	GPT-4o	34.0 \pm 2.9	35.5 \pm 2.5
	Mistral	21.5 \pm 0.3	23.2 \pm 0.4
multi-modal	Llama-3	23.5 \pm 2.5	21.7 \pm 2.0
	GPT-4v	34.2 \pm 1.6	32.0 \pm 1.5
	GPT-4o	52.3\pm0.8	36.5 \pm 0.4
	+CoT	46.0 \pm 2.5	39.5\pm1.1
	InternVL-Chat-V1-5	26.3 \pm 1.6	24.8 \pm 1.3
	LLaVA-NeXT-34B	26.2 \pm 1.1	24.7 \pm 3.2
	Humans	92.0 \pm 4.3	77.8 \pm 6.3

Table 3: Performance of different text-only and multi-modal LLMs on our tasks.

an accuracy rate higher than 90%. The PHYSICO-ASSOCIATIVE tasks present greater challenges and subjectivity because the annotations are personalized based on the annotators’ individual perspectives and experiences. Despite these challenges, humans can still achieve a notable average accuracy of 77.8% in solving these tasks.

We conducted a detailed investigation into human performance on a subset of PHYSICO-ASSOCIATIVE. Participants were asked to annotate instances where they believed none of the four candidate answers adequately explained the inputs. The results revealed a 10.4% rate of disagreement. On these disagreed-upon examples, human accuracy was 33.3%, explaining a major factor for the human performance decline.

Conclusion Our study demonstrates that humans can perform the PHYSICO tasks quite well.

RQ 3: Can LLMs understand concepts in the abstract representations of the matrix format?

A straightforward solution for our PHYSICO is to represent the grid-formatted examples as matrices. By representing the matrices with a token sequence, they can be integrated into an instruction prompt for text-based LLMs, following existing prompting methods for ARC tasks (Acquaviva et al., 2022; Xu et al., 2023; Wang et al., 2023, 2024). We use the prompt shown in Figure 7 to query the answers from the evaluated LLMs.

Results The top (text-only) section of Table 3 presents the results. All the LLMs perform poorly on the two sets of our PHYSICO. Notably, GPT-3.5, Mistral, and Llama-3 failed to show significant improvement over random performance.

Conclusion Comparing the human performance in RQ 2 to the best-performing LLMs reveals a huge gap. While these tasks are simple or trivial for humans, LLMs face substantial challenges, indicating a lack of deep understanding.

When comparing LLMs’ performance on low-level natural language tasks in RQ 1 to high-level abstract pattern understanding tasks, we observe significant declines. This highlights the presence of the *stochastic parrot* phenomenon in LLMs. Our dataset also *quantifies the severity of this phenomenon*. For example, while GPT-3.5 performs on par with GPT-4 on the low-level tasks, it nearly drops to random guessing on our high-level tasks, revealing its tendency to act as a stochastic parrot with the physical concepts in our dataset.

RQ 4: Can multimodal LLMs perform well on our tasks with visual input representations?

Next, we explore whether multi-modal LLMs can effectively solve our tasks when the input examples are presented as visual images rather than matrices like in RQ 3. We use the prompt in Figure 8 to query the answers from evaluated LLMs.

Results The bottom (multi-modal) section of Table 3 shows the results. Consistent with the observations in RQ 3, a significant gap between the performance of LLMs and humans exists.

Notably, the recently introduced GPT-4o outperforms all other LLMs on PHYSICO-CORE by 10% with visual inputs but lags behind GPT-4 on matrix inputs. This discrepancy may be due to GPT-4o’s training on images that directly illustrate physical concepts, making it more adept at solving problems like in Figure 1. However, this advantage does not extend to the more abstract problems in PHYSICO-ASSOCIATIVE that require further knowledge application skills, highlighting the LLMs’ lack of deep understanding even with multi-modal training.⁵

Finally, given that LLMs can generate high-quality descriptions of the concepts (see RQ 1), we adopt a chain-of-thought (Wei et al., 2022) approach. It first asks the LLMs to describe each choice and then makes predictions. The results in Table 3 (+CoT) show limited improvement or performance drop, further confirming the presence of the stochastic parrot phenomenon.

RQ 5: Is PHYSICO challenging mainly due to LLMs’ unfamiliarity with grid representations?

⁵This suggests that the *stochastic parrot phenomenon* may also exist in visual understanding, which we left for future work.

One possible reason for the challenges of PHYSICO might be the uncommon nature of the task format (especially the matrix-format inputs) encountered during LLM training, rather than a lack of deep understanding. We aim to disprove this hypothesis through the following experiments:

- (exp 5.1) *ICL on other concepts*. Compare the performance of zero-shot GPT-4 with GPT-4 using in-context learning (ICL) on few-shot examples from concepts other than the assessed one.
- (exp 5.2) *FT on synthetic matrix data*. Compare the open-source LLMs before and after fine-tuning on a large amount of matrix-input data (Appendix D.1)
- (exp 5.3) *FT on the ARC task*. Compare the open-source LLMs before and after fine-tuning on the original ARC (Chollet, 2019) task, which ensures that all inputs from the PHYSICO-ASSOCIATIVE examples have been seen during model training.

Results and Conclusion Despite that both the ICL and SFT approaches make LLMs more familiar with matrix-format inputs, neither approach significantly improves the results as shown in Table 4. It confirms that merely familiarizing LLMs with grid-format inputs does not enhance their performance on our tasks.

RQ 6: How much can LLMs benefit from training on labeled data?

Many tasks that challenge LLMs can see significant performance boosts through ICL or SFT on labeled data (Hessel et al., 2023; Yu et al., 2023; Berglund et al., 2023). When such improvements are observed, it suggests that LLMs inherently possess the necessary skills to excel in their tasks, needing only minimal training effort.

In this study, we demonstrate that this is not the case for our tasks, where light-weight training on labeled data does not improve LLM performance for our tasks. Given the current lack of large-scale training data for our purpose, we conduct an extreme case study: models learn from the same concepts in PHYSICO-CORE and are tested on the same concepts in PHYSICO-ASSOCIATIVE. To this end, we select the instances that consists of at least two choices that exist in the PHYSICO-CORE, leaving 80 examples. We conduct the following experiments on this subset to answer RQ 6:

- (exp 6.1) *ICL on the same concepts*. Compare the zero-shot GPT-4 and GPT-4 with ICL on examples for the same concepts from PHYSICO-CORE. Specifically, for each instance, we sample

Models	CORE	ASSOCIATIVE
GPT-4	41.3 \pm 1.3	39.0 \pm 0.6
w/ ICL-3-shot	39.5 \pm 1.6	36.2 \pm 1.7
w/ ICL-9-shot	32.8 \pm 1.0	39.0 \pm 1.6
Mistral	21.5 \pm 0.3	23.2 \pm 0.4
w/ FT on syn-tasks	20.9 \pm 0.7	22.5 \pm 0.5
w/ FT on ARC	20.9 \pm 0.8	25.5 \pm 0.9
Llama-3	23.5 \pm 2.5	21.7 \pm 2.0
w/ FT on syn-tasks	23.0 \pm 1.1	23.2 \pm 2.7
w/ FT on ARC	22.2 \pm 1.6	22.4 \pm 1.2

Table 4: Performance of LLMs with in-context learning or fine-tuning on grid-format data.

GPT-4	42.9 \pm 2.4	Llama-3	22.1 \pm 2.8
+ ICL on CORE	40.0 \pm 1.0	+ SFT on CORE	20.9 \pm 2.7

Table 5: Accuracy on PHYSICO-ASSOCIATIVE’s subset that has overlapped choices with PHYSICO-CORE.

9 examples from PHYSICO-CORE with their labels among the choices of the instance.

- (*exp 6.2*) *SFT on the CORE set.* Compare the open-source LLMs before and after fine-tuning on labeled data from PHYSICO-CORE.

Results Table 5 shows that ICL and SFT on the labeled examples of the same concepts even hurt the performance. This is likely because the models are overfitted to the “clean” examples from the PHYSICO-CORE. The difficulty of generalization *within the same concepts* indicates the challenges of our tasks to the supervised fine-tuning paradigm.

Conclusion Together with the results for RQ 5 and RQ 6, it suggests that the low performance of LLMs is not likely to be improved from prompting techniques alone. There exists intrinsic inefficiency in the pre-training of LLMs, which results in the lack of necessary skills for deep understanding.

7 Related Work

Stochastic Parrots on LLMs The pioneer study by (Bender and Koller, 2020) questioned the understanding ability of large models; and Bender et al. (2021) first introduced the terminology of stochastic parrot. The concept of stochastic parrot has received great attention, leading to a surge of studies on this topic. According to Google Scholar, the term “stochastic parrot” appears in the titles of dozens of papers from diverse research fields (Borji, 2023; Li, 2023; Duan et al., 2024; Henrique et al., 2023). However, although the concept of stochastic parrots in LLMs is widely accepted and recognized, to the best of our knowledge, there is a lack

of quantitative experiments to precisely verify this viewpoint. This gap directly motivates our work.

Abstract Reasoning Challenge Abstract reasoning challenge (ARC) aims to examine the inductive reasoning ability in a few-shot scenario (Chollet, 2019): a system is required to generate the output grid for an input grid given a set of input-output examples. ARC has been used as a remarkable testbed to measure the intelligence of LLMs. Recently, many research efforts have been made on improving the performance of LLMs on ARC benchmark (Tan and Motani, 2023; Wang et al., 2023; Xu et al., 2023; Mirchandani et al., 2023; Wang et al., 2024; Huang et al., 2024).

We draw inspiration from ARC by utilizing input-output grids as abstract representations in our task design. However, our task is significantly different from the ARC-style work — our high-level understanding task focuses on comprehending the transformation rules from inputs to outputs and relating them to physical concepts, and is designed to assess the stochastic parrot phenomenon in LLMs.

Challenging Tasks towards LLMs’ Understanding

Extensive recent efforts have been made on designing tasks that challenge the understanding abilities of LLMs (Chakrabarty et al., 2022; Tong et al., 2024; Shapira et al., 2023; Hessel et al., 2023; Donadel et al., 2024; Li et al., 2024). For example, Hessel et al. (2023) proposed a humor understanding task, revealing a large performance gap between LLMs and humans.

As a by-product, our PHYSICO challenges the understanding capabilities of LLMs, relating it to the above studies. However, we make primary contribution to provide an quantitative experiment to verify stochastic parrots in LLMs via controllably paired low-level and high-level tasks.

8 Conclusion

We introduce PHYSICO, a novel task to assess machines’ understanding of physical concepts at different levels. Our experiments reveal that: 1) LLMs lag significantly behind humans on PHYSICO, indicating a lack of deep understanding of the covered concepts; 2) LLMs exhibit the stochastic parrot phenomenon, as they excel at low-level remembering tasks but struggle with high-level understanding tasks; 3) LLMs’ poor performance stems from its intrinsic deficiencies, as neither in-context learning nor fine-tuning improves their results.

686 Limitations

687 Our proposed dataset creation approach is gener-
688 ally extendable but has some limitations. First,
689 the creation of PHYSICO-ASSOCIATIVE relies
690 on the existing ARC dataset, which, despite
691 ongoing additions ([https://arc-editor.
692 lab42.global/playground](https://arc-editor.lab42.global/playground)), offers a lim-
693 ited number of usable instances. Second, the grid
694 world’s representation strength is limited, making it
695 challenging to illustrate many phenomena. Finally,
696 the current PHYSICO-ASSOCIATIVE includes sub-
697 jective cases that are difficult to obtain (see RQ 2
698 for details), introducing noises. In the future, we
699 will continue expanding our tasks along the path
700 of PHYSICO-CORE, by broadening the scope of
701 concepts and their corresponding examples and ex-
702 ploring more ways to represent examples in deeper
703 understanding levels.

704 References

- 705 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
706 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
707 Diogo Almeida, Janko Altenschmidt, Sam Altman,
708 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
709 *arXiv preprint arXiv:2303.08774*.
- 710 Sam Acquaviva, Yewen Pu, Marta Kryven, Theodoros
711 Sechopoulos, Catherine Wong, Gabrielle Ecanow,
712 Maxwell Nye, Michael Tessler, and Josh Tenenbaum.
713 2022. Communicating natural programs to humans
714 and machines. *Advances in Neural Information Pro-
715 cessing Systems*, 35:3731–3743.
- 716 Patricia Armstrong. 2010. Bloom’s taxonomy. *Vander-
717 bilt University Center for Teaching*, pages 1–3.
- 718 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
719 automatic metric for mt evaluation with improved cor-
720 relation with human judgments. In *Proceedings of
721 the acl workshop on intrinsic and extrinsic evaluation
722 measures for machine translation and/or summariza-
723 tion*, pages 65–72.
- 724 Emily M Bender, Timnit Gebru, Angelina McMillan-
725 Major, and Shmargaret Shmitchell. 2021. On the
726 dangers of stochastic parrots: Can language models
727 be too big? In *Proceedings of the 2021 ACM confer-
728 ence on fairness, accountability, and transparency*,
729 pages 610–623.
- 730 Emily M Bender and Alexander Koller. 2020. Climbing
731 towards nlu: On meaning, form, and understanding
732 in the age of data. In *Proceedings of the 58th an-
733 nual meeting of the association for computational
734 linguistics*, pages 5185–5198.
- 735 Lukas Berglund, Meg Tong, Max Kaufmann, Mikita
736 Balesni, Asa Cooper Stickland, Tomasz Korbak, and

- Owain Evans. 2023. The reversal curse: Llms trained
on " a is b" fail to learn" b is a". *arXiv preprint
arXiv:2309.12288*. 737
738
739
- Paul Black and Dylan Wiliam. 1998a. Assessment and
classroom learning. *Assessment in Education: prin-
ciples, policy & practice*, 5(1):7–74. 740
741
742
- Paul Black and Dylan Wiliam. 1998b. *Inside the black
box: Raising standards through classroom assess-
ment*. Granada Learning. 743
744
745
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli,
Russ Altman, Simran Arora, Sydney von Arx,
Michael S Bernstein, Jeannette Bohg, Antoine Bosse-
lut, Emma Brunskill, et al. 2021. On the opportuni-
ties and risks of foundation models. *arXiv preprint
arXiv:2108.07258*. 746
747
748
749
750
751
- Ali Borji. 2023. Stochastic parrots or intelligent sys-
tems? a perspective on true depth of understanding in
llms. *A Perspective on True Depth of Understanding
in LLMs (July 11, 2023)*. 752
753
754
755
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, et al. 2020. Language models are few-shot
learners. *Advances in neural information processing
systems*, 33:1877–1901. 756
757
758
759
760
761
- Tuhin Chakrabarty, Arkadiy Saakyan, Debanjan Ghosh,
and Smaranda Muresan. 2022. **FLUTE: Figurative
language understanding through textual explanations**.
In *Proceedings of the 2022 Conference on Empiri-
cal Methods in Natural Language Processing*, pages
7139–7159, Abu Dhabi, United Arab Emirates. As-
sociation for Computational Linguistics. 762
763
764
765
766
767
768
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,
Cunxiang Wang, Yidong Wang, et al. 2024. A sur-
vey on evaluation of large language models. *ACM
Transactions on Intelligent Systems and Technology*,
15(3):1–45. 769
770
771
772
773
774
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye,
Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi
Hu, Jiapeng Luo, Zheng Ma, et al. 2024. How far
are we to gpt-4v? closing the gap to commercial
multimodal models with open-source suites. *arXiv
preprint arXiv:2404.16821*. 775
776
777
778
779
780
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su,
Guo Chen, Sen Xing, Muyan Zhong, Qinglong
Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo,
Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl:
Scaling up vision foundation models and aligning
for generic visual-linguistic tasks. *arXiv preprint
arXiv:2312.14238*. 781
782
783
784
785
786
787
- François Chollet. 2019. On the measure of intelligence.
arXiv preprint arXiv:1911.01547. 788
789

790	Dorottya Demszky, Diyi Yang, David S Yeager, Christopher J Bryan, Margaret Clapper, Susannah Chandhok, Johannes C Eichstaedt, Cameron Hecht, Jeremy Jamieson, Meghann Johnson, et al. 2023. Using large language models in psychology. <i>Nature Reviews Psychology</i> , 2(11):688–701.	Jiangnan Li, Qiuqing Wang, Liyan Xu, Wenjie Pang, Mo Yu, Zheng Lin, Weiping Wang, and Jie Zhou. 2024. Previously on the stories: Recap snippet identification for story reading. <i>arXiv preprint arXiv:2402.07271</i> .	847
791			848
792			849
793			850
794			851
795			
796	Denis Donadel, Francesco Marchiori, Luca Pajola, and Mauro Conti. 2024. Can llms understand computer networks? towards a virtual system administrator. <i>arXiv preprint arXiv:2404.12689</i> .	Zihao Li. 2023. The dark side of chatgpt: legal and ethical challenges from stochastic parrots and hallucination. <i>arXiv preprint arXiv:2304.14347</i> .	852
797			853
798			854
799			
800	Haonan Duan, Adam Dziedzic, Nicolas Papernot, and Franziska Boenisch. 2024. Flocks of stochastic parrots: Differentially private prompt learning for large language models. <i>Advances in Neural Information Processing Systems</i> , 36.	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.	855
801			856
802			857
803		Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. In <i>NeurIPS</i> .	858
804			859
805	Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. 2022. A survey of embodied ai: From simulators to research tasks. <i>IEEE Transactions on Emerging Topics in Computational Intelligence</i> , 6(2):230–244.	MetaAI. 2024. Introducing meta llama 3: The most capable openly available llm to date .	860
806			861
807			
808		Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. <i>ACM Computing Surveys</i> , 56(2):1–40.	862
809			863
810	Wynne Harlen and Mary James. 1997. Assessment and learning: differences and relationships between formative and summative assessment. <i>Assessment in education: Principles, policy & practice</i> , 4(3):365–379.		864
811			865
812			866
813			867
814		Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, Danny Driess, Montserrat Gonzalez Arenas, Kanishka Rao, Dorsa Sadigh, and Andy Zeng. 2023. Large language models as general pattern machines . <i>ArXiv preprint, abs/2307.04721</i> .	868
815	Da Silva Gameiro Henrique, Andrei Kucharavy, and Rachid Guerraoui. 2023. Stochastic parrots looking for stochastic parrots: LLMs are easy to fine-tune and hard to detect with other LLMs. <i>arXiv preprint arXiv:2304.08968</i> .		869
816			870
817			871
818			872
819			
820	Jack Hessel, Ana Marasovic, Jena D. Hwang, Lillian Lee, Jeff Da, Rowan Zellers, Robert Mankoff, and Yejin Choi. 2023. Do androids laugh at electric sheep? humor “understanding” benchmarks from the new yorker caption contest . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 688–714, Toronto, Canada. Association for Computational Linguistics.	Melanie Mitchell and David C Krakauer. 2023. The debate over understanding in ai’s large language models. <i>Proceedings of the National Academy of Sciences</i> , 120(13):e2215907120.	873
821			874
822			875
823			876
824		Kishore Papineni, Salim Roukos, Todd Ward, and Weijng Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	877
825			878
826			879
827			880
828			881
829	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. Lora: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	Di Peng, Liubin Zheng, Dan Liu, Cheng Han, Xin Wang, Yan Yang, Li Song, Miaoying Zhao, Yanfeng Wei, Jiayi Li, et al. 2024. Large-language models facilitate discovery of the molecular signatures regulating sleep and activity. <i>Nature Communications</i> , 15(1):3685.	882
830			883
831			884
832			885
833			886
834	Di Huang, Ziyuan Nan, Xing Hu, Pengwei Jin, Shaohui Peng, Yuanbo Wen, Rui Zhang, Zidong Du, Qi Guo, Yewen Pu, et al. 2024. Anpl: Towards natural programming with interactive decomposition. <i>Advances in Neural Information Processing Systems</i> , 36.	Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. 2019. Habitat: A platform for embodied ai research. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 9339–9347.	887
835			888
836			889
837			890
838			891
839	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	Natalie Shapira, Guy Zwirn, and Yoav Goldberg. 2023. How well do large language models perform on faux pas tests? In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 10438–10451, Toronto, Canada. Association for Computational Linguistics.	892
840			893
841			894
842			895
843			896
844	David R Krathwohl. 2002. A revision of bloom’s taxonomy: An overview. <i>Theory into practice</i> , 41(4):212–218.		897
845			898
846			899

900 John Chong Min Tan and Mehul Motani. 2023. [Large](#)
901 [language model \(llm\) as a system of multiple expert](#)
902 [agents: An approach to solve the abstraction and](#)
903 [reasoning corpus \(arc\) challenge](#). *ArXiv preprint*,
904 [abs/2310.05146](#).

905 Gemini Team, Rohan Anil, Sebastian Borgeaud,
906 Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
907 Radu Soricut, Johan Schalkwyk, Andrew M Dai,
908 Anja Hauth, et al. 2023. Gemini: a family of
909 highly capable multimodal models. *arXiv preprint*
910 [arXiv:2312.11805](#).

911 Xiaoyu Tong, Rochelle Choenni, Martha Lewis, and
912 Ekaterina Shutova. 2024. Metaphor understand-
913 ing challenge dataset for llms. *arXiv preprint*
914 [arXiv:2403.11810](#).

915 Ruocheng Wang, Eric Zelikman, Gabriel Poesia, Yewen
916 Pu, Nick Haber, and Noah D Goodman. 2023. [Hy-](#)
917 [pothesis search: Inductive reasoning with language](#)
918 [models](#). *ArXiv preprint*, [abs/2309.05660](#).

919 Yile Wang, Sijie Cheng, Zixin Sun, Peng Li, and Yang
920 Liu. 2024. [Speak it out: Solving symbol-related](#)
921 [problems with symbol-to-language conversion for](#)
922 [language models](#). *ArXiv preprint*, [abs/2401.11725](#).

923 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
924 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
925 et al. 2022. Chain-of-thought prompting elicits rea-
926 soning in large language models. *Advances in neural*
927 *information processing systems*, 35:24824–24837.

928 Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui
929 Wang, Zichao Yang, and Zhiting Hu. 2023. Language
930 models meet world models: Embodied experiences
931 enhance language models. *Advances in neural infor-*
932 *mation processing systems*, 36.

933 Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott
934 Sanner, and Elias B Khalil. 2023. [Llms and the](#)
935 [abstraction and reasoning corpus: Successes, failures,](#)
936 [and the importance of object-based representations](#).
937 *ArXiv preprint*, [abs/2305.18354](#).

938 Mo Yu, Jiangnan Li, Shunyu Yao, Wenjie Pang, Xi-
939 aochen Zhou, Zhou Xiao, Fandong Meng, and Jie
940 Zhou. 2023. Personality understanding of fictional
941 characters during book reading. In *Proceedings*
942 *of the 61st Annual Meeting of the Association for*
943 *Computational Linguistics (Volume 1: Long Papers)*,
944 pages 14784–14802.

A Details of the Included Concepts in our PHYSICO

Concepts in PHYSICO-CORE The concepts in PHYSICO-CORE are basic physical concepts that we manually design problems for. The set covers 27 concepts as follows:

reference frame	12	gravity	10
reflection	10	refraction	10
light imaging	10	communicating vessels	10
cut	10	laser	10
surface tension	10	move	10
buoyancy	10	acceleration	10
inertia	10	electricity	10
repulsive force	8	wave	8
lever	6	optical filters	6
compression	4	diffuse reflection of light	4
wave interference	4	diffusion	4
vortex	4	expansion	4
nuclear fission	2	nuclear fusion	2
diffraction of waves	2		

Table 6: Concepts and their corresponding number of instances in PHYSICO-CORE.

All Concepts in PHYSICO The following table summarized all the concepts from both PHYSICO-CORE and PHYSICO-ASSOCIATIVE:

laser	30	mirror	30
wave	21	reference frame	20
gravity	19	move	18
reflection	15	zoom in	15
compression	14	magnet	14
expansion	13	explosion	11
refraction	10	light imaging	10
communicating vessels	10	cut	10
surface tension	10	buoyancy	10
acceleration	10	inertia	10
electricity	10	rotation	10
repulsive force	8	diffusion	8
optical filters	7	water ripples	7
long exposure	7	lever	6
wave interference	5	vortex	5
wetting	5	diffuse reflection of light	4
nuclear fission	3	nuclear fusion	3
zoom out	3	diffraction of waves	2
projection	2	polarization of light	1
chemical bond	1	squeeze	1
lumination	1	vacuum	1

Table 7: Concepts and their corresponding number of instances in PHYSICO-CORE.

B Details of Analysis Methods in RQ 1

B.1 Masking of Textual Descriptions

This experiment follows the setting in the ‘‘Physical Concept Selection Subtask’’ in section 3.2.1. The definitions of the corresponding phenomena were extracted from Wikipedia as well as generated by GPT-3.5 and GPT-4. To maintain consistency,

the terms representing concepts were masked as [PHENOMENON] while relevant terms are masked as [MASK]. For instance, ‘‘interference’’ which corresponds to the phenomenon ‘‘wave interference’’ was masked as [PHENOMENON]. In contrast, ‘‘Newton’s first law of motion’’ which corresponds to the phenomenon ‘‘inertia’’ was masked as [MASK].

An example of the masked description can be found in Figure 6.

B.2 Prompts Used for Description Generation and Classification

```
[SYSTEM]
You are an expert in physics. Your task is
↪ to provide a comprehensive definition of
↪ a given physical concept or phenomenon,
↪ with the key properties or key examples
↪ of the concept included.

[USER]
Please provide me with the definition of
↪ the physical concept "{ CONCEPT }",
↪ with the key properties or key examples
↪ included.
```

Figure 4: The prompt template used for generating descriptions of physical concepts (denoted as the variable **CONCEPT**) in RQ 2.

```
[SYSTEM]
You will be playing a game:
You are given a definition of a physical
↪ phenomenon, where the names of the
↪ phenomenon are masked.
Your task is to guess which phenomenon the
↪ definition refers to.
Please select the most close answer from
↪ the provided options.

[USER]
Here is a definition of a physical
↪ phenomenon, where the names of the
↪ phenomenon are masked:

[Definition]
{{ MASKED DESCRIPTION }}

Please guess which phenomenon the
↪ definition refers to. You should choose
↪ your answer from the following options:
↪ {{ CANDIDATE ANSWERS }}

Your response should end with your choice
↪ of answer.
```

Figure 5: The prompt template used for guessing the referred physical concept from four candidates (denoted as the variable **CANDIDATE ANSWERS**) from the natural language descriptions (denoted as the variable **MASKED DESCRIPTION**) in RQ 2.

```

[PHENOMENON] is a fundamental concept in physics that describes the resistance of any physical
↪ object to a change in its state of motion. This concept is a central part of [MASK], often
↪ referred to as the law of [PHENOMENON]. According to this law, an object at rest will stay at
↪ rest, and an object in motion will continue to move at a constant velocity, unless acted upon by
↪ a net external force. Here are the key properties and examples of [PHENOMENON]:

### Key Properties:
1. **Dependence on Mass**: The [PHENOMENON] of an object is directly proportional to its mass. The
↪ greater the mass of an object, the greater its [PHENOMENON], and hence, the more force it
↪ requires to change its state of motion.
2. **Resistance to Acceleration**: [PHENOMENON] is essentially the resistance of an object to any
↪ change in its velocity, which includes changes in the speed or direction of the object's motion.
3. **Universal Applicability**: [PHENOMENON] applies to all objects with mass, whether they are
↪ microscopic or astronomical in scale.
4. **Independence from External Factors**: The [PHENOMENON] of an object is inherent and does not
↪ depend on external conditions such as the environment, temperature, or pressure.

### Key Examples:
1. **A Parked Car**: A parked car will not move unless a force is applied to it. Once moving, it
↪ will continue to move at a constant speed in a straight line unless forces like friction or
↪ brakes are applied to change its state.
2. **Astronauts and Objects in Space**: In the vacuum of space, where there is little to no
↪ external force, an astronaut or any other object will continue moving in the same direction and
↪ at the same speed until acted upon by another force. This is an example of [PHENOMENON] in a
↪ microgravity environment.
3. **Seatbelts in Vehicles**: When a car suddenly stops, the passengers inside tend to lurch
↪ forward. This is due to the [PHENOMENON] of their bodies; their bodies were in motion and tend to
↪ remain in motion despite the car stopping. Seatbelts provide the necessary force to counteract
↪ this [PHENOMENON] and keep the passengers safe.
4. **Tablecloth Trick**: A classic example demonstrating [PHENOMENON] is the tablecloth trick,
↪ where a quick pull of the tablecloth can leave dishes undisturbed on a table. The [PHENOMENON] of
↪ the dishes (their tendency to resist changes in motion) allows them to remain relatively still
↪ while the tablecloth is quickly pulled from under them.

Understanding [PHENOMENON] is crucial for analyzing the motion of objects in various physical
↪ contexts, from everyday life to complex scientific scenarios. It is a cornerstone in the study of
↪ dynamics and plays a critical role in engineering, automotive safety, aerospace technology, and
↪ many other fields.

```

Figure 6: An example of our masked description for the concept inertia.

C Details of the Methods Used in RQ 3 and RQ 4

We use the prompt template in Figure 7 for experiments on text-only LLMs (RQ 3); and the template in Figure 8 for multi-modal LLMs (RQ 4).

D Details of Supervised LLM Fine-Tuning in RQ 5 and RQ 6

D.1 Construction of Synthetic Training Data in RQ 5

We investigate whether fine-tuning LLMs on matrix property-related questions could improve their performances on our tasks. Specifically, we generate 3000 extra input-output grid pairs calculate the size, transpose, and locations of the subgrid’s corner elements for these matrices as ground truths. Furthermore, since correctly recognizing the location of the subgrid may contribute more to finish the Move and Copy tasks compared to other properties, we create additional ground truths only with the gold locations of the subgrid’s corner elements.

D.2 Training Details for RQ 5 and RQ 6

For all the fine-tuning experiments, we use LoRA (Hu et al., 2021). We fine-tune each model for 3 epochs with a batch size of 4 on a single machine with 8 A100 GPUs. The dimension of LoRA’s attention layer is set to 64, and the α and dropout rates are set to 16 and 0.1, respectively. The learning rate and weight decay are set to $2e-4$ and 0.001, respectively. The hyperparameters are selected according to the development performance on the synthetic matrix data in Appendix D.1.

```

[SYSTEM]
You will be playing a game:
You are given several examples. Each example consists of an ``input grid'' and an ``output grid'' of
↔ numbers from 0-9, where each number corresponds to a color.
Your task is to try to find the common patterns from the examples and abstract the meanings of the
↔ patterns in the physical or mathematics world.
Based on the recognized meaning, please select the most close description of the common pattern
↔ from the provided options.

[USER]
Let's play a game where you are transforming an input grid of numbers into an output grid of numbers.
↔

The numbers represent different colors:
0 = black
1 = blue
2 = red
3 = green
4 = yellow
5 = gray
6 = magenta
7 = orange
8 = cyan
9 = brown

Here are examples of input grids and its corresponding output grids:

Example input grid:
{{ INPUT GRID1 }}

Example output grid:
{{ OUTPUT GRID1 }}

Example input grid:
{{ INPUT GRID2 }}

Example output grid:
{{ OUTPUT GRID2 }}

Example input grid:
{{ INPUT GRID3 }}

Example output grid:
{{ OUTPUT GRID3 }}

Please first try to find the common patterns from the input-output pairs, then answer the following
↔ question:

What meanings in the physical or mathematics world can be abstracted from the patterns? Please
↔ choose your answer from the following options:
{{ CANDIDATE ANSWERS }}

Your response should end with your choice of answer.

```

Figure 7: The prompt template used in RQ 3. The pair of an **INPUT GRID** and an **OUTPUT GRID** consists of one example of a physical phenomenon in matrix format.

```

{{ UPLOADED IMAGE }}
[USER]
In the given image, there are two columns of matrices with elements represented by different colors.
↔
The left column represents the inputs, and the right column represents the corresponding outputs.
For each row in the image, the output is derived from the input using the same transformation rule,
which corresponds to a real-world physical concept.

Your task is to identify the physical concept demonstrated in this image from the following options:
↔

{{ CANDIDATE ANSWERS }}

Please select and provide the correct option that matches the transformation shown in the image.
Your response should end with your choice of answer.

```

Figure 8: The prompt template used in RQ 4. **UPLOADED IMAGE** is an image consists of three or more examples like in Figure 2.