

SCI PREDICT: CAN LLMs PREDICT THE OUTCOMES OF SCIENTIFIC EXPERIMENTS IN NATURAL SCIENCES?

Anonymous authors

Paper under double-blind review

ABSTRACT

Accelerating scientific discovery requires the identification of which experiments would yield the best outcomes before committing resources to costly physical validation. While existing benchmarks evaluate LLMs on scientific knowledge and reasoning, their ability to predict experimental outcomes—a task where AI could significantly exceed human capabilities—remains largely unexplored. We introduce SciPredict, a benchmark comprising 405 tasks derived from recent empirical studies in 33 specialized sub-fields of physics, biology, and chemistry. SciPredict addresses two critical questions: (a) *can LLMs predict the outcome of scientific experiments with sufficient accuracy?* and (b) *can such predictions be reliably used in the scientific research process?* Evaluations reveal fundamental limitations on both fronts. Model accuracies are 14-26% and human expert performance is $\approx 20\%$. Although some frontier models exceed human performance model accuracy is still far below what would enable reliable experimental guidance. Even within the limited performance, models fail to distinguish reliable predictions from unreliable ones, achieving only $\approx 20\%$ accuracy regardless of their confidence or whether they judge outcomes as predictable without physical experimentation. Human experts, in contrast, demonstrate strong calibration: their accuracy increases from $\approx 5\%$ to $\approx 80\%$ as they deem outcomes more predictable without conducting the experiment. SciPredict establishes a rigorous framework demonstrating that superhuman performance in experimental science requires not just better predictions, but better awareness of prediction reliability. For reproducibility all our data and code are provided at <https://anonymous.4open.science/r/SciPredict-AI01>.

1 INTRODUCTION

Reasoning deeply about the expected outcome of experiments before running them is central to efficient scientific progress Platt (1964). Researchers routinely make such predictions, deciding which hypotheses to test and parameter regimes to pursue under resource constraints. A system that could reliably predict the experimental results would reshape the scientific process, accelerating discovery by filtering out suboptimal directions, identifying gaps in current frameworks, and suggesting much needed empirical investigations. LLMs appear well-suited for this task (as illustrated in Appendix A.1 Fig. 11), as they encode vast scientific knowledge Taylor et al. (2022), can reason about complex systems, and demonstrate strong performance on scientific QA tasks Wang et al. (2025).

Due to the lack of comprehensive benchmarks, the progress toward improving the ability of LLMs to predict the outcomes of practical experiments has been slow. Among benchmarks that explore the use of LLMs to aid the scientific research, most focus on areas such as literature review and paper composition/drafting Li et al. (2025); Xu et al. (2025); Laurent et al. (2024), reproducing methods and computational simulation results Starace et al. (2025); Zhao et al. (2025); Yan et al. (2025); Lin et al. (2025); Ali-Dib & Menou (2024); Shojaee et al. (2025); Xia et al. (2025), or generating hypotheses for scientific experiments Yang et al. (2025); Ke et al. (2025); Abdel-Rehim et al. (2025). A comprehensive discussion on related work is provided in Appendix A.2.

To address this gap, we introduce SciPredict, a benchmark designed to evaluate the capabilities of LLMs in predicting the outcomes of empirical experiments in natural sciences. We extract tasks from recently published empirical studies, from post-March 31, 2025, postdating the cutoff dates of frontier models. SciPredict comprised of 405 experimental prediction tasks, spanning 33 specialized

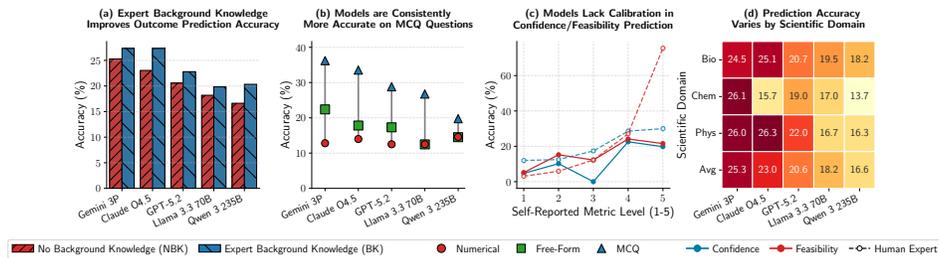
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Figure 1: **Key findings of SciPredict.** Frontier models exhibit fundamental gaps in accuracy and calibration robustness in scientific experiment outcome prediction. We highlight four key failure modes using a representative subset of SOTA models: Claude O4.5 (Claude Opus 4.5), OpenAI GPT-5.2, Gemini 3P (Gemini 3 Pro), Llama 3.3 (Meta Llama 3.3 70B), and Qwen 3 235B. (a) Providing expert-curated background knowledge (BK) consistently boosts performance over No Background Knowledge (NBK). (b) Accuracy generally degrades when moving from multiple-choice questions (MCQ) to questions requiring free-form answers to Numerical value questions. (c) Unlike Human Experts (dashed lines), models show poor calibration in SciPredict tasks; the accuracy of the models’ answers to tasks do *not* correlate with their self-reported Confidence and perceived task prediction Feasibility. (d) Model performance varies across different domains. The Avg field shown represents the weighted average of scores according to the number of questions per domain.

sub-fields: 9 under physics, 10 under chemistry, and 14 under biology. For each task, domain-expert human annotators extract relevant information, including experimental setups, measurements taken by the research team, empirical results, etc. from the target publications along with the relevant background knowledge from prior literature. Prediction questions come in three possible formats of multiple-choice (MCQ), free-format (FF), or numerical value (NUM) depending on the task. This variation allows us to effectively capture the different aspects of models’ capabilities in scientific reasoning. For free-format questions SciPredict includes 1-10 expert-written rubrics used to judge the accuracy of provided predictions. For MCQ and NUM questions, respectively, the correct choice(s) and acceptable numerical ranges are provided as the ground-truth labels. Each task underwent a multi-stage expert review process. The curation process overall costed \$336k and 7,380 human expert hours, reflecting the difficulty of constructing a high-quality benchmark for experimental outcome prediction.

Our findings show that SOTA LLMs achieve prediction accuracy between 14%-26% while human experts achieve $\approx 20\%$. Although exceeding human performance, these accuracy levels remain insufficient for reliable experimental planning. In practice, the *reliability* of the outcome prediction process is crucial because researchers want to invest their limited resources in sufficiently compelling experimental directions. To account for this, we require the models and human experts to provide prediction *feasibility* scores along with their predictions, measuring whether the targeted outcomes are perceived to be reliably predictable given the contextual information (e.g., experimental setup, background information), without physically conducting the experiments. Models show poor calibration of such scores with their measured prediction accuracy: their accuracy does *not* meaningfully improve with higher self-reported feasibility scores. Human experts, on the other hand, demonstrate *strong* calibration of their prediction accuracy and their rated prediction feasibility scores (increase in accuracy from $\approx 5\%$ to $\approx 80\%$ as rated feasibility rises).

To understand what types of prior scientific knowledge aids accurate outcome predictions, we used different variations of *background knowledge* in our evaluations. While expert-curated background knowledge (mainly extracted by experts from prior studies cited in the target publication) improved accuracy by $\approx 3\%$ on average (1.2 – 5.8% depending on the model), the models’ self-generated background knowledge often resulted in accuracy degradation. Interestingly, even combining such self-generated background knowledge items with the expert-curated knowledge still yielded under-performance in most cases. We note that this pattern reveals a critical limitation: models struggle to identify what background information and prior scientific knowledge would be helpful for task outcome predictions, often introducing misleading assumptions or irrelevant details in their self-generated background knowledge that degrades accuracy. Fig. 1 summarizes some of our primary findings.

Our key contributions are summarized as follows:

- We introduce SciPredict, the first benchmark for evaluating LLMs in experimental outcome prediction tasks in natural sciences (biology, chemistry, physics). This dataset is comprised of 405 expert-curated tasks with three prediction question types (multiple-choice, free-form, and numerical) directly derived from empirical studies published after March 31, 2025, ensuring no data leakage from the model pre-training data.
- We conduct a comprehensive evaluation of 15 SOTA LLMs and human experts, analyzing the accuracy and reliability (confidence, difficulty, feasibility). We analyze the effectiveness of 4 types of relevant background knowledge being provided in context for effective predictions (expert-curated, self-generated, filtered, combined).
- We identify a critical calibration gap: unlike human experts who demonstrate strong calibration of their confidence/difficulty/feasibility ratings with their prediction accuracy, LLMs mostly do not show such meaningful correlations, making their deployment in real-world scientific experimentation pipelines untrustworthy.
- We demonstrate that the models benefit from expert-curated background knowledge provided in context for predictions, while they fail to generate such background knowledge autonomously. We also reveal primary causes for models prediction failures are due to factual and logical reasoning flaws rather than misunderstanding the task.

2 SCIPREDICT CURATION

SciPredict consists of 405 prediction tasks derived from empirical studies published after March 2025 across physics, biology, and chemistry. The construction process balances competing requirements. Questions must be challenging enough to distinguish model capabilities yet tractable enough that expert-curated background knowledge could plausibly aid prediction. Experimental setups must be described with sufficient precision for informed reasoning without simply revealing the answer. Ground truth outcomes must be objectively verifiable while accounting for the inherent variability in empirical measurements.

Design Principles. Our focus domains are experimentally rich and empirical validation is central to knowledge creation. To evaluate scientific reasoning, we use three question formats—multiple-choice (MCQ), free-form, and numerical—to cover discrete, explanatory, and quantitative prediction. Domain-selection criteria, rubric designs, and evaluation specifics are detailed in Appendix B.2. Example data are given in Appendix C.1.

Data Collection. We recruited domain experts across biology, physics, and chemistry, representing diverse educational and geographic backgrounds (details in Appendix B.1). Experts selected papers published after March 31, 2025 ensuring tasks represent genuine predictive reasoning challenges. They extracted experimental setups, measurements, prediction questions, and ground-truth outcomes, and curated background knowledge relevant for informed reasoning. Full details of selection criteria and extraction methods are provided in Appendix B.3.

Quality Control. All tasks underwent rigorous multi-stage review. Initial screening removed ambiguous, simulated, theoretical, or outdated (pre-March 2025) tasks. Two rounds of domain experts verified the clarity and completeness of experimental details, background knowledge relevance, ground truth clarity, and task difficulty. Reviewers additionally ensured that MCQ distractors represented

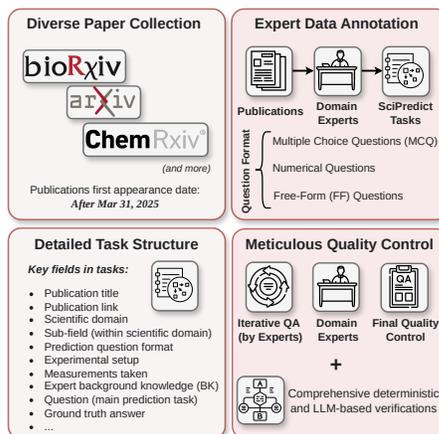


Figure 2: **SciPredict curation pipeline.** The benchmark construction involves four integrated stages: (Top-Left) **Data Collection** from preprint repositories ensuring a post-March 2025 cutoff to prevent data leakage; (Top-Right) **Expert Annotation** where domain specialists convert raw papers into MCQ, numerical, and free-form prediction tasks; (Bottom-Left) **Task Structure** enforcement, ensuring every sample includes granular fields such as experimental setup, measurements, and expert-curated background knowledge; and (Bottom-Right) **Quality Control/Assurance**, combining iterative expert review with deterministic and LLM-based verifications.

scientifically sound but incorrect alternatives, comprehensive yet flexible free-form evaluation rubrics, and realistic numerical precision ranges. See Appendix B.5 for full reviewer guidelines and checks.

Data Diversity. SciPredict spans 33 sub-fields of physics, biology, and chemistry. Task distribution across these domains are provided in Appendix B.7 Tab. 2. Tasks systematically vary in complexity, from single-step causal reasoning to complex multi-hop inference requiring advanced expertise. Background knowledge items principles range from undergraduate-level to very specialized (expertise held by active researchers). Task distribution ensures sufficient representation across domains (physics 25%, biology 50%, chemistry 25%) and question formats (MCQ 40%, free-form 32%, numerical 28%). See Appendix B.6 for additional details.

Human baseline. In addition to the experts who constructed the benchmark, we recruited a separate group of domain experts to provide a human baseline performance. Each expert answered benchmark questions, provided reasoning for their answers, and reliability scores. Mirroring the LLM evaluation, experts completed two rounds: first without, then with provided background knowledge (further details are provided in Appendix B.4). To ensure a high quality human baseline, we match experts to tasks based on their expertise. Additional details on the mapping and example human baseline responses are given in Appendix B.7 Tab. 1 and C.2.

3 EVALUATION SETUP AND METRICS

Our dataset \mathcal{D} contains 3 subsets for multiple-choice questions (\mathcal{D}_{MCQ}), free-form questions (\mathcal{D}_{FF}), and numerical questions (\mathcal{D}_{NUM}). For each task i , evaluated models $m \in \mathcal{M}$ provide a prediction $\hat{y}_i^{(m)}$ and 3 reliability assessments.

3.1 ACCURACY METRICS

Multiple-choice (MCQ). Each question $i \in \mathcal{D}_{\text{MCQ}}$ presents 3-4 options with ground truth answer $g_i \subseteq \{A, B, C, D\}$ provided by domain expert annotators (g_i is comprised of more than one choice in 12.35% of the MCQ tasks only). Accuracy is the proportion of questions answered correctly: $\text{Acc}_{\text{MCQ}}^{(m)} = \frac{1}{|\mathcal{D}_{\text{MCQ}}|} \sum_{i \in \mathcal{D}_{\text{MCQ}}} \mathbb{1}[\hat{y}_i^{(m)} = g_i]$.

Free-form (FF). Each question $i \in \mathcal{D}_{\text{FF}}$ has a reference answer y_i and an expert-written evaluation rubric. We employ an LLM judge J_θ with a fixed prompt to assess whether the model’s response $\hat{y}_i^{(m)}$ demonstrates correct scientific reasoning, as evaluated against the provided rubric: $s_i = J_\theta(\hat{y}_i^{(m)}, y_i) \in \{0, 1\}$, $\text{Acc}_{\text{FF}}^{(m)} = \frac{1}{|\mathcal{D}_{\text{FF}}|} \sum_{i \in \mathcal{D}_{\text{FF}}} s_i$.

Numerical value (NUM). For each question $i \in \mathcal{D}_{\text{NUM}}$, domain experts specify an acceptable range $[L_i, U_i]$ accounting for measurement precision and experimental variability. Accuracy reflects whether predictions fall within this scientifically reasonable interval: $\text{Acc}_{\text{NUM}}^{(m)} = \frac{1}{|\mathcal{D}_{\text{NUM}}|} \sum_{i \in \mathcal{D}_{\text{NUM}}} \mathbb{1}[L_i \leq \hat{y}_i^{(m)} \leq U_i]$. This metric captures whether the model’s quantitative prediction is sufficiently accurate for experimental planning, rather than demanding exact numerical matches.

3.2 RELIABILITY CALIBRATION

Reliable deployment in experimental science requires not only accurate predictions but also the ability to distinguish trustworthy predictions from unreliable ones. We assess reliability through three complementary measures.

- **Confidence.** Models report confidence $\hat{c}_i^{(m)} \in \{1, \dots, 5\}$ regarding their prediction’s correctness. If well-calibrated, this metric is expected to *positively* correlate with the prediction accuracy ($\hat{c} \uparrow$ correlates with $\text{Acc} \uparrow$)
- **Difficulty.** Models’ perceived task prediction hardness $\hat{z}_i^{(m)} \in \{1, \dots, 5\}$ given the provided context. Difficulty assesses the self-awareness of models regarding their own prediction limitations. If well-calibrated, this metric is expected to *negatively* correlate with the prediction accuracy ($\hat{z} \uparrow$ correlates with $\text{Acc} \downarrow$).

- **Feasibility.** Models assess if an outcome can be predicted via reasoning without running the practical experiment ($\hat{f}_i^{(m)} \in \{1, \dots, 5\}$). If well-calibrated, this metric is expected to *positively* correlate with the prediction accuracy ($\hat{f} \uparrow$ correlates with Acc \uparrow).

3.3 EXPERIMENTAL CONDITIONS

To determine the information requirements for accurate predictions, we systematically vary the **context** provided to the model. Each task’s BK in SciPredict is comprised of multiple atomic knowledge bullet points. We evaluate under five conditions:

- **No Background Knowledge (NBK).** The context contains only the experimental setup, measurements, and the prediction question. This assesses whether the model’s *parametric knowledge* is sufficient for prediction.
- **Background Knowledge (BK).** The context additionally includes expert-curated BK. This measures the performance gain when relevant, high-quality background information is explicitly surfaced in the context.
- **Self-generated Background (SBK).** The model is prompted to generate its own BK before predicting. This assesses the model’s ability to autonomously identify and articulate the necessary scientific context.
- **Self-generated + Annotator Background (SABK).** The context includes both the expert-curated (BK) and self-generated background knowledge (SBK). This assesses whether combining such information sources provides additive benefits or introduces noise/interference.
- **Filtered Background Knowledge (FBK).** For each model, the context includes expert BK *minus* the facts the model already *knows*. We convert the BK items into questions and remove any BK items from the final prediction context where the models is able to answer the corresponding questions. This isolates whether stating known information in context improves prediction even when that information is theoretically accessible from parameters.

3.4 EVALUATION PROTOCOL AND ROBUSTNESS

Free-form predictions were evaluated by Gemini 3-Pro against expert rubrics. We validated the robustness of such evaluation pipeline by replicating evaluations using GPT-5.2 as well, where we found *no statistically significant* differences in accuracy scores. We also replicated predictions using various decoding strategies (temperature settings from 0.0 to 1.0, top-p sampling with $p \in \{0.9, 0.95, 1.0\}$). Performance variations remained statistically insignificant. Reported accuracy metrics represent means and with error bars indicate one standard deviation within 3 trials.

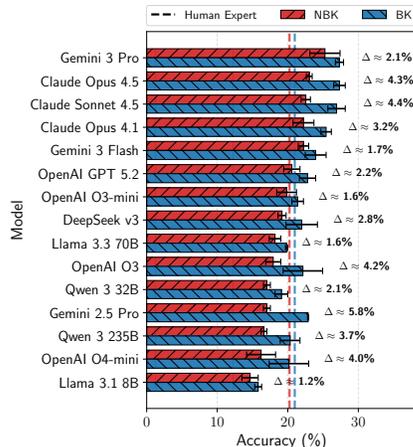


Figure 3: **Accuracy with and without background knowledge.** Accuracy (%) of each evaluated model under two input conditions: (a) **W/o background knowledge**: the model receives only the experimental setup, measurements, and the question; (b) **W/ background knowledge**: the same information as previous case with the addition of expert-curated background knowledge.

4 MAIN RESULTS

We assess 15 SOTA LLMs analyzing if frontier models can predict experimental outcomes with sufficient accuracy and reliability for scientific use. All the prompts used in evaluations are given in Appendix E.

Finding #1: Human performance is close to the average model performance.

We emphasize that expert human baseline performance serves as a calibration reference point, *not* an upper bound; the models can exceed human prediction capabilities by integrating vast cross-domain knowledge and reasoning power. Human baseline ($\approx 20\%$ accuracy; Fig. 3) reflects the inherent difficulty of predicting novel experimental outcomes without real-world scientific experimentation or validation.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Finding #2: Providing curated background knowledge consistently improves the outcome prediction accuracy.

A key factor in answering the questions correctly, is access to relevant background knowledge. We test two conditions, model performance with and without background knowledge. As shown in Fig. 3, providing background knowledge improves accuracy $Acc^{(m)}$ across all models $m \in \mathcal{M}$, though the size of the increase varies by model. On average, BK improves accuracy by $\approx 3\%$. Curated background knowledge helps narrow the space of plausible outcomes. It is noted that confidence scores $\hat{c}^{(m)}$ remain roughly the same across NBK and BK, suggesting that background information primarily benefits correctness rather than improving confidence.

Finding #3: Across nearly all models, accuracy is higher with the full annotator background than with the filtered version, implying that including knowledge the models already know still boosts performance.

Fig. 4 shows that restating known facts in the input context enhances model performance, even when those facts are *not* strictly missing from the models’ parametric knowledge. By filtering the background knowledge—removing any expert background knowledge bullet points that the models already demonstrate knowledge of; see §Sec. 3.3—the x-axis approximates performance when the context contains only the “unknown” background knowledge. Most models fall in the upper triangle (above the $y = x$ line), illustrating accuracy $Acc^{(m)}$ is higher when the full curated background is provided, including facts the model demonstrably knows. Additional results are given in Appendix D Tab. 3.

Finding #4: Models cannot reliably generate useful background knowledge: self-generated background usually reduces accuracy. When combined with expert-curated background knowledge it rarely improves performance.

We evaluate settings where models self-generate background knowledge (SBK) and then answer, as well as a combined condition that appends this self-generated context to expert-curated background (SABK). Fig. 6 shows that, in contrast to the clear gains from expert-curated background knowledge, self-generated background is unreliable and often counterproductive: for most models, SBK lowers accuracy compared to providing no background at all, implying that the generated content is frequently irrelevant or misleading and can steer predictions away from the correct experimental

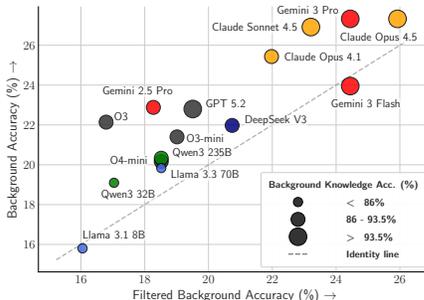


Figure 4: **Restating known facts in context enhances performance.** We present a scatter plot comparing model accuracy given full expert-curated background (BK, y-axis) versus a filtered version of such background (FBK, x-axis) in context (see definitions under §Sec. 3.3). Each point represents a model, colored by family, with marker size indicating the percentage of the tasks’ background knowledge bullet points the model already knew ($> 70\%$ for all). Most models lie above the dashed identity line ($y = x$), showing that explicitly restating already known required background knowledge in the context yields higher accuracy than relying on originally learned knowledge alone.

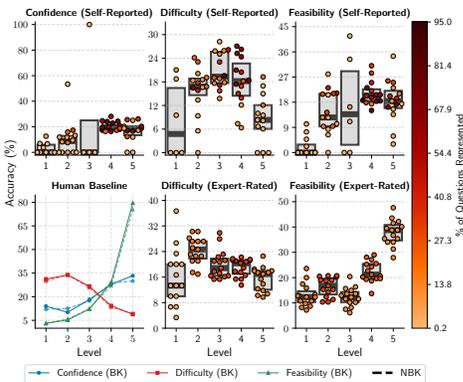


Figure 5: **Models are poorly calibrated in self-reported confidence, difficulty, and feasibility, whereas human calibration correlated with accuracy.** The top row plots empirical accuracy against model-provided confidence/difficulty/feasibility metrics; the expected trends (accuracy rising with confidence/feasibility and falling with difficulty) are weak and non-monotonic. The bottom row shows the prediction accuracy of models and humans plotted against the human calibration metrics. The bottom-left subplot shows human accuracy with significant positive correlation with rated confidence and feasibility levels, while showing a significant negative correlation with the rated difficulty as expected. The bottom-middle and bottom-right subplots reveal that model accuracy recovers the expected trends when plotted against *human-rated* difficulty and feasibility levels, confirming that human judgment provides a superior signal for task outcome predictability compared to models’ self-reports. Circle colors correspond to the percentage of the number of total questions assigned to that calibration level (x-axis) by each model (top row) or human expert (bottom row). Each circle corresponds to a distinct model.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

outcome. Moreover, supplementing expert-curated background knowledge with self-generated background (SABK) typically fails to yield consistent improvements, indicating that models struggle not only to generate helpful knowledge, but also to avoid introducing distracting or harmful information when additional context is available. Additional results are given in Appendix D Tab. 3.

Finding #5: Self-assessed confidence/difficulty/feasibility by models are not aligned with accuracy, indicating calibration gaps. Humans, in contrast, show strong calibration of their rated confidence/difficulty/feasibility scores with their accuracy.

Fig. 5 demonstrates if models $m \in \mathcal{M}$ can reliably anticipate their own prediction errors by comparing accuracy $\text{Acc}^{(m)}$ to self-reported confidence $\hat{c}^{(m)}$, difficulty $\hat{z}^{(m)}$, and feasibility $\hat{f}^{(m)}$ ratings. While informative self-assessments would be expected to imply that $\text{Acc}^{(m)} \uparrow$ as $\hat{c}^{(m)} \uparrow$, $\hat{f}^{(m)} \uparrow$, and $\hat{z}^{(m)} \downarrow$, our results show that these relationships are very weak and non-monotonic for the models, which indicates substantial *miscalibration*. Conversely, human ratings of confidence, difficulty, and feasibility closely track accuracy in the expected directions (strong positive or negative correlation). Interestingly, as the bottom-left subfigure of Fig. 5 shows, models systematically achieve higher accuracy ($\text{Acc}_i^{(m)} \uparrow$) on tasks where human experts rate as more feasible ($f_i^{(m)} \uparrow$) and less difficult ($z_i^{(m)} \downarrow$), demonstrating that the human experts can much more reliably capture the predictability of evaluated tasks. Additional results are given in Appendix D Tab. 5 and Tab. 6.

Finding #6: Model failures are primarily driven by factual and extraction errors (avg. 80.14%) and logical reasoning flaws (avg. 87.42%) rather than comprehension issues, frequently manifesting as information fabrication and false certainty.

We classify model failures into 16 error types and 5 categories (Fig. 7) using an LLM judge. Results show failures concentrate in Factual and Extraction errors (avg. 80.1%) and Logical and Reasoning flaws (avg. 87.4%). Prevalent fine-grained errors include Factual Contradiction (avg. 52.3%) and Information Fabrication (avg. 54.0%), indicating that models frequently fail to incorporate relevant experimental information and basic scientific facts when making predictions. Deficiencies in Scientific Rigor (avg. 47.9%) appear primarily as False Certainty (avg. 43.6%), where models express establish probabilistic or rough predictions as certain facts. Models also on average fail to acknowledge their limitations in providing predictions in about 19.4% of the tasks. Basic comprehension/scope (avg. 10.0%) and formatting/mechanical (avg. <0.6%) errors remain rare, confirming models understand the tasks but lack reasoning capabilities for effective predictions. These patterns persist when only considering tasks which human ex-

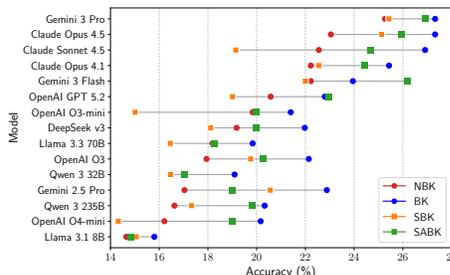


Figure 6: **Human vs self-generated background knowledge.** Evaluated accuracy (%) for the models under the four prediction conditions defined in §Sec. 3.3. *BK* generally yields the highest accuracy, while *SBK* frequently degrades accuracy relative to *NBK*, indicating that models fail to reliably generate useful predictive context. Furthermore, *SABK* rarely improves upon *BK*, suggesting that adding synthetic information likely introduces noise or misleading cues even when the correct expert background information is available in context.

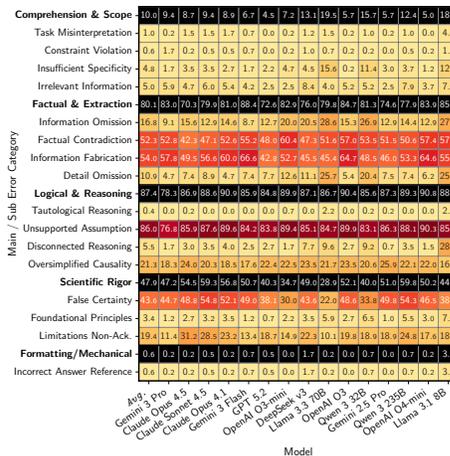


Figure 7: **Analysis of model errors.** We errors into a hierarchical taxonomy spanning five top-level (in black background) categories and 16 specific error types. The heatmap shows the percentage of incorrect responses containing each error type for each evaluated model. Error categories (as defined in Tab. 8) progress from surface-level issues (e.g., Comprehension & Scope) to deeper reasoning failures (e.g., Logical & Reasoning flaws) to fundamental scientific deficiencies (Scientific Rigor flaws). Models can exhibit multiple error types simultaneously, so accumulative percentage scores within top-level categories may exceed 100%. SciPredict tasks contribute to top-level category percentages if flagged with at least one underlying error type. Fig. 13 shows the same chart only for tasks with human rated *feasibility* $\in \{4, 5\}$.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

perts rate as feasible ($\hat{f} \in \{4, 5\}$) as shown in Fig. 13. Tab. 8 provides detailed definitions for the error categories.

Finding #7: MCQs are substantially easier than free-form and numerical value tasks.

As shown in Fig. 8 we find that model accuracy is highly sensitive to answer format, with multiple-choice questions substantially easier than open-ended generation and especially numerical prediction. This gap is not merely a matter of “MCQs being easier because the correct option is visible,” but appears to reflect a broader dependence on recognition over generation: MCQs let models compare candidates and pick the closest match, while free-form and numerical formats require constructing a specific claim/value and committing to it. To isolate format from content, we convert MCQs into matched free-form prompts (MCQ→FF) and re-run evaluation. The resulting drop, visible across essentially all model families, shows that simply removing the provided options degrades accuracy even when the underlying experimental scenario is unchanged. This suggests that headline MCQ accuracy $\text{Acc}_{\text{MCQ}}^{(m)}$ can overestimate how reliably a model would perform in realistic scientific workflows, where predictions are typically produced in open form $\text{Acc}_{\text{FF}}^{(m)}$ (and often as quantities). Finally, the steepness of the MCQ→free-form drop varies by model, implying meaningful differences in robustness to output constraints. Additional results are given in Appendix D Tab. 4.

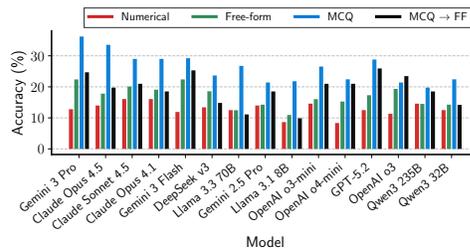


Figure 8: **Accuracy is highly sensitive to the question format.** Even when the underlying prediction tasks are identical, performance is generally the highest for multiple-choice questions (MCQs), lower for free-form (FF), and lowest for numerical (NUM) predictions. Crucially, when MCQs are rewritten as free-form (MCQ → FF) questions to remove the answer options, accuracy drops. This shows that a model’s reported performance depends heavily on how the prediction is requested.

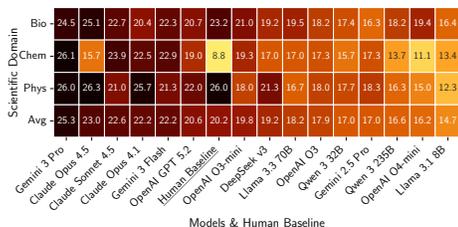


Figure 9: **Domain specific accuracy.** Heatmap of model accuracy (%) on benchmark questions, broken down by scientific domain (Biology, Physics, Chemistry). Results are provided for the evaluated models and human baseline. Overall, frontier models achieve the highest accuracies, but performance varies by domain; Chemistry tends to be the most challenging subset, and several models (including the human baseline) exhibit performance degradation on Chemistry relative to Biology/Physics. The Avg shown represents the weighted average respective to number of questions per domain.

anticipate empirical outcomes from experimental descriptions. Although the overall association with HLE is positive, the dispersion around the trendline is substantial: models with similar HLE text-only accuracy can differ by several points on NBK accuracy. This residual structure is informative: some models overperform relative to what their HLE score would predict (e.g., DeepSeek v3 achieves comparatively strong NBK accuracy despite very low HLE, and Claude Sonnet 4.5 / Claude Opus 4.1 sit above the fitted line), while others underperform given their HLE level (e.g., Gemini 2.5 Pro, OpenAI o3, and GPT-5.2 fall below the line). These deviations suggest that, beyond general text-only reasoning, strong results on SciPredict also depend on scientific priors and experimental intuition: identifying which intervention details are causally relevant, mapping measurements to plausible mechanisms, and remaining robust when background context is withheld in the NBK setting.

Finding #8: Performance varies by scientific domain, with Chemistry typically being the most challenging.

Fig. 9 shows that Chemistry has the lowest accuracy on average compared to Biology and Physics. This domain gap is particularly visible for the human baseline. Accuracy gains of models are not uniform across domains, indicating that scaling or general instruction-following ability does not fully translate into robust empirical reasoning in Chemistry. This pattern suggests that our benchmark is sensitive to domain-specific experimental knowledge and intuitions. Additional results are given in Appendix D Tab. 3.

Finding #9: Performance on this benchmark has a strong correlation with performance on HLE benchmark.

Fig. 10 helps disentangle how much performance on SciPredict (NBK) reflects broad hard-reasoning capability versus a more task-specific ability to

5 DISCUSSION AND CONCLUSIONS

Our work reveals fundamental gaps between current LLM capabilities and the requirements for reliable experimental guidance. While frontier models achieve 14-26% accuracy, comparable to human expert baselines around 20%, performance remains insufficient for guiding resource-intensive experimental decisions. Models exhibit severe miscalibration: unlike human experts whose accuracy ranges from ~5% on infeasible questions to ~80% on feasible ones, models maintain uniform ~20% performance regardless of self-reported confidence or feasibility. Expert-curated background knowledge provides modest gains (~3%), but models cannot autonomously identify or generate helpful context. These findings demonstrate that achieving superhuman scientific assistance requires not merely better predictions, but systems that accurately assess their own reliability.

Limitations. While SciPredict establishes a rigorous framework for evaluating experimental outcome prediction, some limitations constrain the scope and generalizability. The benchmark focuses on 3 natural science domains, excluding engineering and computational fields where prediction tasks may exhibit different characteristics. Our temporal cutoff (March 2025) ensures data freshness but limits historical coverage, and the 405-question scale, though substantial, may not capture the full diversity of experimental paradigms within each subdomain. The reliance on expert-curated background knowledge, while ensuring quality, introduces potential biases in what information is deemed relevant.

Future work. The path toward AI systems that meaningfully accelerate scientific discovery extends beyond improving prediction accuracy on static benchmarks. Integrating models with active experimentation frameworks would enable systems to propose experiments, observe outcomes, and iteratively refine hypotheses, transforming prediction from a one-shot task into a dialogue between theory and empirical validation. Developing methods for cross-domain knowledge transfer could allow models to recognize when principles from one field apply to another, mimicking how expert scientists draw analogies across disciplines.

IMPACT STATEMENT

This work introduces SciPredict, a benchmark for evaluating AI systems' ability to predict scientific experimental outcomes, a capability that could significantly accelerate scientific discovery by helping researchers prioritize experiments and allocate resources more efficiently. The potential benefits of achieving reliable experimental outcome prediction are substantial: reduced costs and time in experimental research, faster iteration cycles in scientific discovery, and more efficient allocation of limited research resources, particularly in domains where experiments are costly or time-consuming. However, our work also highlights critical risks that must be addressed before deployment, most notably severe miscalibration where models express high confidence on incorrect predictions and cannot distinguish reliable forecasts from unreliable ones. Deploying such systems prematurely could lead researchers to pursue unproductive experimental directions, waste valuable resources, or—in high-stakes domains like medicine or materials science—make decisions with serious real-world consequences based on unreliable AI predictions. We emphasize that overcoming the calibration gap we identify is as important as improving raw accuracy, and our benchmark provides the measurement framework necessary to track progress toward AI scientific assistants that are not only capable but also trustworthy and appropriately cautious about their own limitations.

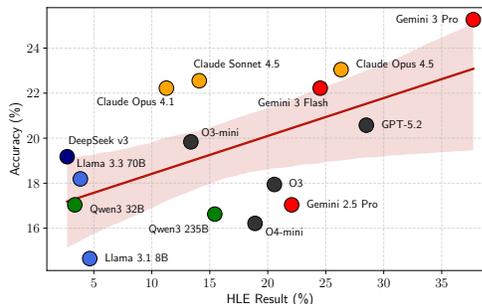


Figure 10: **Model accuracy on SciPredict correlates with performance on the HLE benchmark.**

Benchmark performance correlates with general hard-reasoning performance. This is a scatter plot of each evaluated model's accuracy on SciPredict in the no-background-knowledge (NBK) setting (y axis) versus its HLE text-only accuracy (x axis). The solid line shows a linear fit and the shaded region indicates the corresponding confidence bands. Overall, SciPredict NBK accuracy exhibits a moderate positive correlation with HLE performance (Pearson $r \approx 0.46$), suggesting that broader reasoning capability explains some-but not all-variance in empirical outcome prediction.

REFERENCES

- 486
487
488 Abbi Abdel-Rehim, Hector Zenil, Oghenejokpeme Orhobor, Marie Fisher, Ross J Collins, Elizabeth
489 Bourne, Gareth W Fearnley, Emma Tate, Holly X Smith, Larisa N Soldatova, et al. Scientific
490 hypothesis generation by large language models: laboratory validation in breast cancer treatment.
491 *Journal of the Royal Society Interface*, 22(227):20240674, 2025.
- 492 Mohamad Ali-Dib and Kristen Menou. Physics simulation capabilities of llms. *Physica Scripta*, 99
493 (11):116003, 2024.
- 494 Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhui Chen, and William Yang Wang. Knowledge
495 of knowledge: Exploring known-unknowns uncertainty with large language models. In *Findings of*
496 *the Association for Computational Linguistics: ACL 2024*, pp. 6416–6432, 2024.
- 498 Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela,
499 Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, et al.
500 Healthbench: Evaluating large language models towards improved human health. *arXiv preprint*
501 *arXiv:2505.08775*, 2025.
- 502 Casey O Barkan, Sid Black, and Oliver Sourbut. Do large language models know what they are
503 capable of? *arXiv preprint arXiv:2512.24661*, 2025.
- 505 Dennis Bersenev, Ayako Yachie-Kinoshita, and Sucheendra K Palaniappan. Replicating a high-impact
506 scientific publication using systems of large language models. *bioRxiv*, pp. 2024–04, 2024.
- 508 Daniel Brunnsåker, Alexander H Gower, Prajakta Naval, Erik Y Bjurström, Filip Kronström,
509 Ievgeniia A Tiukova, and Ross D King. Self-driven biological discovery through automated
510 hypothesis generation and experimental validation. *bioRxiv*, pp. 2025–06, 2025.
- 511 Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio
512 Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Mađry.
513 Mle-bench: Evaluating machine learning agents on machine learning engineering, 2025. URL
514 <https://arxiv.org/abs/2410.07095>.
- 515 Hui Chen, Miao Xiong, Yujie Lu, Wei Han, Ailin Deng, Yufei He, Jiaying Wu, Yibo Li, Yue Liu, and
516 Bryan Hooi. Mlr-bench: Evaluating ai agents on open-ended machine learning research, 2025.
517 URL <https://arxiv.org/abs/2505.19955>.
- 519 Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema
520 Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. Finqa: A dataset
521 of numerical reasoning over financial data, 2022. URL <https://arxiv.org/abs/2109.00122>.
- 522 Ziyang Cui, Ning Li, and Huaikang Zhou. Can ai replace human subjects? a large-scale replication of
523 psychological experiments with llms. *A Large-Scale Replication of Psychological Experiments*
524 *with LLMs (August 25, 2024)*, 2024.
- 526 Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A
527 comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.
- 528 Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex
529 Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry
530 Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai
531 Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin Tobia,
532 Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam
533 Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi,
534 Tom Zur, Varun Iyer, and Zehua Li. Legalbench: A collaboratively built benchmark for measuring
535 legal reasoning in large language models, 2023. URL <https://arxiv.org/abs/2308.11462>.
- 536 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks.
537 In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on*
538 *Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1321–1330.
539 PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/guo17a.html>.

- 540 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
541 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
542
- 543 Tianyu Hua, Harper Hua, Violet Xiang, Benjamin Klieger, Sang T Truong, Weixin Liang, Fan-Yun
544 Sun, and Nick Haber. Researchcodebench: Benchmarking llms on implementing novel machine
545 learning research code. *arXiv preprint arXiv:2506.02314*, 2025.
- 546 Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents
547 on machine learning experimentation, 2024. URL <https://arxiv.org/abs/2310.03302>.
548
- 549 Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. How can we know when language
550 models know? on the calibration of language models for question answering. *Transactions of the*
551 *Association for Computational Linguistics*, 9:962–977, 2021.
- 552 Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and
553 Yuxiang Wu. Aide: Ai-driven exploration in the space of code, 2025. URL <https://arxiv.org/abs/2502.13138>.
554
- 555 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W. Cohen, and Xinghua Lu. Pubmedqa: A dataset
556 for biomedical research question answering, 2019. URL <https://arxiv.org/abs/1909.06146>.
557
- 558 Lennart Justen. Llms outperform experts on challenging biology benchmarks. *arXiv preprint*
559 *arXiv:2505.06108*, 2025.
560
- 561 Yujing Ke, Kevin George, Kathan Pandya, David Blumenthal, Maximilian Sprang, Gerrit Großmann,
562 Sebastian Vollmer, and David Antony Selby. Biodisco: Multi-agent hypothesis generation with
563 dual-mode evidence, iterative feedback and temporal evaluation. *arXiv preprint arXiv:2508.01285*,
564 2025.
- 565 Patrick Tser Jern Kon, Jiachen Liu, Xinyi Zhu, Qiuyi Ding, Jingjia Peng, Jiarong Xing, Yibo Huang,
566 Yiming Qiu, Jayanth Srinivasa, Myungjin Lee, et al. Exp-bench: Can ai conduct ai research
567 experiments? *arXiv preprint arXiv:2505.24785*, 2025.
568
- 569 Jon M. Laurent, Joseph D. Janizek, Michael Ruzo, Michaela M. Hinks, Michael J. Hammerling,
570 Siddharth Narayanan, Manvitha Ponnampati, Andrew D. White, and Samuel G. Rodrigues. Lab-
571 bench: Measuring capabilities of language models for biology research, 2024. URL <https://arxiv.org/abs/2407.10362>.
572
- 573 Matthew Li, Santiago Torres-Garcia, Shayan Halder, Phani Kuppa, Sean O’Brien, Vasu Sharma,
574 Kevin Zhu, and Sunishchal Dev. Frontierscience bench: Evaluating ai research capabilities in llms.
575 In *Proceedings of the 1st Workshop for Research on Agent Language Models (REALM 2025)*, pp.
576 428–453, 2025.
- 577 Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words.
578 *arXiv preprint arXiv:2205.14334*, 2022.
579
- 580 Zijie Lin, Yiqing Shen, Qilin Cai, He Sun, Jinrui Zhou, and Mingjun Xiao. Autop2c: An llm-based
581 agent framework for code repository generation from multimodal content in academic papers.
582 *arXiv preprint arXiv:2504.20115*, 2025.
- 583 Genglin Liu, Xingyao Wang, Lifan Yuan, Yangyi Chen, and Hao Peng. Examining llms’ uncertainty
584 expression towards questions outside parametric knowledge. *arXiv preprint arXiv:2311.09731*,
585 2023.
586
- 587 Sirui Lu, Zhijing Jin, Terry Jingchen Zhang, Pavel Kos, J Ignacio Cirac, and Bernhard Schölkopf. Can
588 theoretical physics research benefit from language agents? *arXiv preprint arXiv:2506.06214*, 2025.
- 589 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa : A large-
590 scale multi-subject multi-choice dataset for medical domain question answering, 2022. URL
591 <https://arxiv.org/abs/2203.14371>.
592
- 593 John R Platt. Strong inference: Certain systematic methods of scientific thinking may produce much
more rapid progress than others. *science*, 146(3642):347–353, 1964.

- 594 Denitsa Saynova, Kajsa Hansson, Bastiaan Bruinsma, Annika Fredén, and Moa Johansson. Identifying
595 non-replicable social science studies with language models. *arXiv preprint arXiv:2503.10671*,
596 2025.
- 597 Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K Reddy.
598 Llm-sr: Scientific equation discovery via programming with large language models. In *The*
599 *Thirteenth International Conference on Learning Representations*, 2025.
- 600 Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z Ren, and Anirudha Majumdar. A survey on
601 uncertainty quantification of large language models: Taxonomy, open research challenges, and
602 future directions. *ACM Computing Surveys*, 2025.
- 603
604 Nithin Somasekharan, Ling Yue, Yadi Cao, Weichao Li, Patrick Emami, Pochinapeddi Sai Bhargav,
605 Anurag Acharya, Xingyu Xie, and Shaowu Pan. Cfd-llmbench: A benchmark suite for evaluating
606 large language models in computational fluid dynamics. *arXiv preprint arXiv:2509.20374*, 2025.
- 607
608 Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel
609 Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, et al. Paperbench: Evaluating ai’s ability
610 to replicate ai research. In *Forty-second International Conference on Machine Learning*, 2025.
- 611
612 Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony S. Hartshorn, Elvis Saravia,
613 Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science.
614 *ArXiv*, abs/2211.09085, 2022. URL <https://api.semanticscholar.org/CorpusID:253553203>.
- 615
616 George Tsatsaronis, Michael Schroeder, Georgios Paliouras, Yannis Almirantis, Ion Androutopoulos,
617 Eric Gaussier, Patrick Gallinari, Thierry Artieres, Michael R. Alvers, Matthias Zschunke, and
618 Axel-Cyrille Ngonga Ngomo. BioASQ: A challenge on large-scale biomedical semantic indexing
619 and Question Answering. In *Proceedings of AAAI Information Retrieval and Knowledge Discovery*
620 *in Biomedical Text*, 2012.
- 621
622 Miles Wang, Robi Lin, Kat Hu, Joy Jiao, Neil Chowdhury, Ethan Chang, and Tejal Patwardhan. Fron-
623 tierscience: Evaluating AI’s ability to perform expert-level scientific tasks. <https://cdn.openai.com/pdf/2fcd284c-b468-4c21-8ee0-7a783933efcc/frontierscience-paper.pdf>, December
624 2025. Technical report. Accessed: 2026-01-26.
- 625
626 Shijie Xia, Yuhan Sun, and Pengfei Liu. Sr-scientist: Scientific equation discovery with agentic ai.
627 *arXiv preprint arXiv:2510.11661*, 2025.
- 628
629 Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms
630 express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint*
631 *arXiv:2306.13063*, 2023.
- 632
633 Tianze Xu, Pengrui Lu, Lyumanshan Ye, Xiangkun Hu, and Pengfei Liu. Researcherbench: Evaluating
634 deep ai research systems on the frontiers of scientific inquiry. *arXiv preprint arXiv:2507.16280*,
635 2025.
- 636
637 Shuo Yan, Ruochen Li, Ziming Luo, Zimu Wang, Daoyang Li, Liqiang Jing, Kaiyu He, Peilin Wu,
638 George Michalopoulos, Yue Zhang, et al. Lmr-bench: Evaluating llm agent’s ability on reproducing
639 language modeling research. *arXiv preprint arXiv:2506.17335*, 2025.
- 640
641 Daniel Yang, Yao-Hung Hubert Tsai, and Makoto Yamada. On verbalized confidence scores for llms.
642 *arXiv preprint arXiv:2412.14737*, 2024.
- 643
644 Zonglin Yang, Wanhao Liu, Ben Gao, Tong Xie, Yuqiang Li, Wanli Ouyang, Soujanya Poria, Erik
645 Cambria, and Dongzhan Zhou. Large language models for rediscovering unseen chemistry scientific
646 hypotheses. In *2nd AI4Research Workshop: Towards a Knowledge-grounded Scientific Research*
647 *Lifecycle*, 2025.
- 648
649 Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large
650 language models know what they don’t know? *arXiv preprint arXiv:2305.18153*, 2023.
- 651
652 Hanning Zhang, Shizhe Diao, Yong Lin, Yi R Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng
653 Ji, and Tong Zhang. R-tuning: Teaching large language models to refuse unknown questions. *arXiv*
654 *preprint arXiv:2311.09677*, 63:67, 2023.

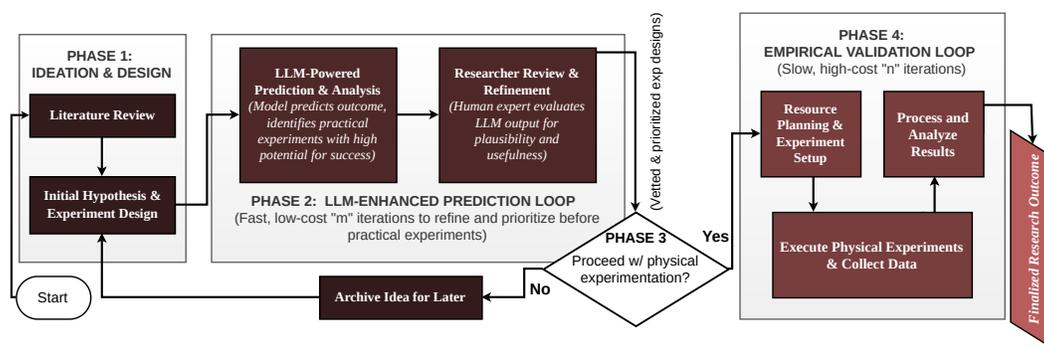
648 Xiao Zhang, Ji Wu, Zhiyang He, Xien Liu, and Ying Su. Medical exam question answering with
649 large-scale reading comprehension, 2018. URL <https://arxiv.org/abs/1802.10279>.

651 Yunxiang Zhang, Muhammad Khalifa, Shitanshu Bhushan, Grant D Murphy, Lajanugen Logeswaran,
652 Jaekyeom Kim, Moontae Lee, Honglak Lee, and Lu Wang. Mlrc-bench: Can language agents solve
653 machine learning research challenges?, 2025. URL <https://arxiv.org/abs/2504.09702>.

654 Xuanle Zhao, Zilin Sang, Yuxuan Li, Qi Shi, Weilun Zhao, Shuo Wang, Duzhen Zhang, Xu Han,
655 Zhiyuan Liu, and Maosong Sun. Autoreproduce: Automatic ai experiment reproduction with paper
656 lineage. *arXiv preprint arXiv:2505.20662*, 2025.

659 A AI FOR SCIENCE

661 A.1 USE OF LLMs IN ACCELERATING SCIENTIFIC PROGRESS



675 **Figure 11: LLM-enhanced efficient scientific research workflow.** The figure illustrates how LLM-
676 powered experimental outcome prediction can be integrated into the scientific research process. Phase
677 1 involves ideation and experimental design through literature review and hypothesis formulation.
678 Phase 2 represents a fast, low-cost prediction loop where LLMs predict experimental outcomes
679 and identify high-potential experiments for physical validation, which researchers then review for
680 plausibility. Based on this evaluation, researchers either proceed to Phase 3 (resource planning and
681 experiment setup) and Phase 4 (empirical validation through physical experimentation), or archive
682 the idea for later consideration. This workflow demonstrates how reliable LLM predictions could
683 accelerate scientific discovery by filtering suboptimal experimental directions before committing to
684 costly empirical validation.

686 A.2 RELATED WORK

687
688 **Expert-level benchmarks in science and professional domains.** Recent studies suggest that
689 LLMs can approach domain experts on selected tasks and in some cases surpass them, while still
690 exhibiting notable gaps in reliability, safety, and grounded reasoning. In scientific computing,
691 end-to-end computational fluid dynamics remains a stringent test of scientific reasoning and code
692 generation, highlighting domain-specific weaknesses that general progress in NLP has not yet closed
693 Somasekharan et al. (2025). In healthcare, LLMs show steady gains in multi-turn evaluations, but
694 important challenges remain for safety-critical decision support Arora et al. (2025). Recent biology
695 evaluations find that frontier LLMs can meet or exceed expert performance on several challenging
696 benchmarks, while also cautioning about benchmark saturation and evaluation errors Justen (2025).
697 Several other benchmarks focus on the evaluation of LLMs in questions from medicine Zhang et al.
698 (2018); Pal et al. (2022); Jin et al. (2019), biomedical research Tsatsaronis et al. (2012), finance Chen
699 et al. (2022), and law Guha et al. (2023). Du et al. (2025) presents a benchmark of 100 PhD-level
700 questions across a broad span of the aforementioned topics. Although these benchmarks require
701 specialized knowledge, they have two primary shortcomings that our work addresses. First, most do
not require the same degree of complex reasoning. Second, they are not situated in the empirical
settings that define our benchmark, which is essential to assess real-world performance.

AI/ML research benchmarks. Recent benchmarks have begun evaluating LLMs on tasks that simulate the AI research cycle itself, extending beyond problem-solving or knowledge recall. Starace et al. (2025); Zhao et al. (2025); Yan et al. (2025); Lin et al. (2025) evaluate LLMs for their ability to reproduce masked or full code repositories and experiment results given existing ML papers. Hua et al. (2025) takes this a step further by evaluating how well LLMs can write experiment code for novel research ideas not seen during training. Huang et al. (2024); Jiang et al. (2025); Chan et al. (2025) evaluate agents on machine learning engineering tasks, assessing their ability to iteratively modify algorithms and improve performance across various datasets and tasks. Li et al. (2025) focuses on research methodology, requiring LLMs to predict masked out methodological details of AI research papers. Xu et al. (2025) evaluates LLM agents' ability to provide technical details, literature review, and open consulting to AI-related questions. Chen et al. (2025); Zhang et al. (2025); Kon et al. (2025) extend evaluation to the entire AI research cycle, asking LLM agents to propose novel ideas or hypotheses, design and execute experiments, and write papers or solutions without a reference. While all of these benchmarks advance the evaluation of LLMs in research-oriented or engineering tasks, they primarily emphasize ideation, writing, or code execution. Our benchmark instead focuses on assessing LLMs' ability to understand and predict empirical scientific outcomes, a skill particularly relevant for research in the physical sciences.

Non-ML scientific research benchmarks. LLMs have also been evaluated for their performance on scientific research tasks outside of AI. For example, Ali-Dib & Menou (2024) assesses LLMs on coding and problem-solving tasks in computational physics. Shojaee et al. (2025) uses LLMs, leveraging their extensive domain knowledge and reliable program synthesis, to infer scientific equations directly from datasets; extending this, Xia et al. (2025) turns LLMs into autonomous scientists that code, evaluate, and iteratively optimize the simulated equations. Similarly, Bersenev et al. (2024) provides LLM agents with written biology papers and evaluates their ability to reproduce the methodology, code, and results. Laurent et al. (2024) tests LLMs on their ability to do literature review and data analysis for biology research questions. While these benchmarks are valuable for evaluating LLMs' abilities in problem-solving, coding, and scientific writing, they do not directly measure an LLM's capacity to predict empirical scientific outcomes.

Work on outcome prediction has so far focused mainly on behavioral and social sciences. Cui et al. (2024) and Saynova et al. (2025) evaluate LLMs on predicting experimental outcomes or reproducibility, but they operate in domains where measurements are often less precise and quantitative. In contrast, our benchmark targets the hard sciences, emphasizing quantitative prediction of empirical results. Lu et al. (2025) provides qualitative analysis of how well LLMs can answer theoretical physics questions using a physics knowledge toolbox, but unlike their position paper, we provide a standardized benchmark for quantitative evaluation.

LLM-driven scientific hypothesis generation. While some benchmarks ask LLMs to generate hypotheses for scientific experiment settings, these works differ from our work in important ways. Yang et al. (2025) provides a benchmark where LLMs have to produce and rank novel hypotheses in chemistry when prompted with background information and a set of hand-picked inspiration facts. Ke et al. (2025) proposes a multi-agent framework that combines language-model reasoning with biomedical knowledge graphs and an automated literature retrieval engine to generate and iteratively refine grounded, novel hypotheses in biomedicine. Abdel-Rehim et al. (2025) examines the applicability of large language models for hypothesis generation, focusing their experiments on breast cancer therapy. Brunnsåker et al. (2025) introduces an LLM-driven approach to automating experimental design that fuses relational learning-generated hypotheses with real-world lab constraints and is deployed on an automated cell and metabolomics platform. While our benchmark also asks LLMs to produce hypotheses in scientific settings, we crucially do not single out inspiration facts, which can heavily influence LLM performance on this task setting.

Confidence evaluation. Confidence can be assessed in two complementary ways: (i) implicit confidence derived from model's output distribution (e.g. logits/probabilities) Jiang et al. (2021); Guo et al. (2017), and (ii) explicit self-reported confidence Lin et al. (2022); Xiong et al. (2023); Yang et al. (2024); Shorinwa et al. (2025). While implicit confidence scores have been proven to provide useful signals for identifying misclassified and out-of-distribution examples Hendrycks & Gimpel (2016), logits are inherently designed to measure the probability of individual tokens rather than full sentences. To solve this, heuristics have been proposed to aggregate token-level scores, but they often

fail to accurately capture the uncertainty over claims themselves Lin et al. (2022). Moreover, implicit methods require access to log-probabilities, which black-box APIs typically do not provide Xiong et al. (2023). For these reasons, we prioritize reporting explicit confidence in this benchmark.

Feasibility evaluation. Self-assessment of feasibility is a classification problem: determining whether a model can successfully complete a given task. Barkan et al. (2025) studies whether LLMs "know what they are capable of" before making an attempt to solve the problem. UnknownBench measures refusal behavior of LLMs using lexical keyword matching over model outputs Liu et al. (2023). Similarly, Yin et al. (2023) determines the LLMs' uncertainty by comparing responses against a set of vague reference sentences via text-similarity scoring. Amayuelas et al. (2024); Zhang et al. (2023) propose datasets that label tasks as feasible or infeasible, enabling a more systematic evaluation of LLM feasibility judgments. Yet, unlike our work, none of these papers evaluate feasibility judgments in the specific setting of empirical outcome prediction, where answering often requires running the underlying experiment.

B ADDITIONAL DATASET DETAILS

B.1 ADDITIONAL DETAILS ABOUT TASK CONTRIBUTORS / HUMAN BASELINE PARTICIPANTS

We provide additional visualizations of the degree, expertise, and country of origin diversity of the experts recruited for benchmark construction and human baseline. Overall, our experts have strong credentials in their respective fields. For the human baseline, we match experts with relevant expertise to task domains and subdomains; see Tab. 1 for more details.

Expert recruitment. To construct our benchmark, we recruited a large cohort of experts in biology, physics, and chemistry. Among them, 54.5% hold a doctoral degree (PhD or equivalent), 34.3% hold a master's degree, and 11.2% hold a bachelor's degree. The experts represent a diverse set of countries, including the United States (14.3%), India (14.3%), United Kingdom (13.6%), Argentina (7.3%), and more. See Fig. 12 for more details.

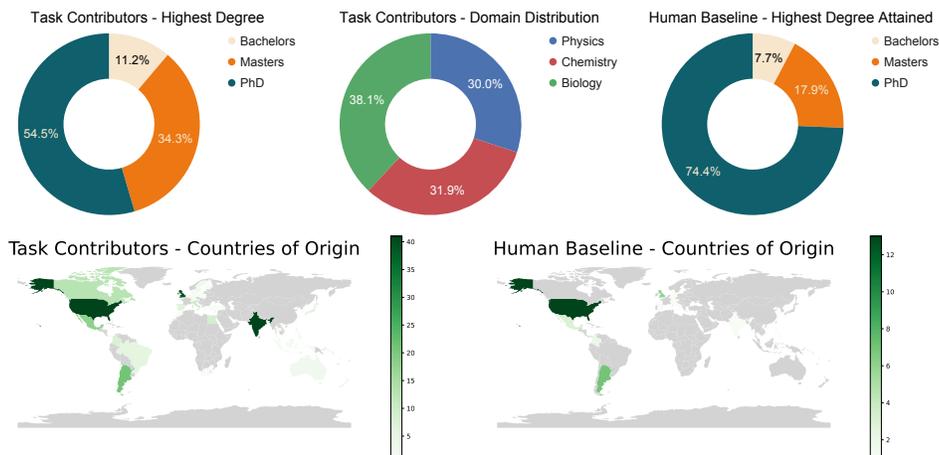


Figure 12: **Diversity of the experts recruited for benchmark construction and human baseline.** *Top left:* A plot of the highest degree distribution of experts recruited for benchmark construction. *Top center:* A plot of the domain expertise of experts recruited for benchmark construction. *Top right:* A plot of the highest degree distribution of experts recruited for human baseline. *Bottom left:* A heatmap of the countries of origin of experts recruited for benchmark construction. *Bottom right:* A heatmap of the countries of origin of experts recruited for human baseline.

B.2 DOMAIN SELECTION AND QUESTION DESIGN DETAILS

Domain Selection Criteria. We selected physics, biology, and chemistry based on three key criteria. First, these domains involve high-stakes applications in engineering, medicine, and materials science,

810 where incorrect predictions can incur significant real-world costs. Second, experimental protocols in
811 these domains are typically well documented, enabling structured extraction of experimental setups,
812 controlled conditions, and measured outcomes. Third, the domains provide sufficient diversity in
813 experimental systems and reasoning styles to evaluate whether models can generalize predictive
814 reasoning across distinct scientific contexts.

815
816 **Question Formats and Evaluation.** Our benchmark includes multiple-choice (MCQ), free-form,
817 and numerical value questions, each with domain-appropriate evaluation procedures. For MCQs,
818 ground truth specifies the correct option or set of correct options. For free-form questions, domain
819 experts design detailed evaluation rubrics that capture the essential scientific reasoning and expected
820 outcomes. For numerical value questions, experts define acceptable answer ranges based on
821 measurement precision and inherent experimental variability, and model predictions are evaluated
822 based on whether they fall within these ranges.

823 B.3 DETAILED EXPERT TASK CURATION

824
825 To prevent data leakage from existing pretraining data, recruited domain experts selected empirical
826 research papers published exclusively after March 31, 2025. Selected studies explicitly avoided purely
827 theoretical analyses or computational simulations, focusing solely on clearly documented empirical
828 experiments. Papers are selected from domain specific open source venues that are widely recognized
829 in the scientific community such as bioRxiv, chemRxiv, arXiv, PubMed Central (PMC), Nature,
830 Science.

831 For each chosen paper, experts explicitly extracted and documented: (1) the domain and specialized
832 subdomain classification, (2) experimental setup details, (3) specific measurements obtained from the
833 experiment, (4) a clear prediction question targeting the experiment’s outcome, and (5) the ground
834 truth answer directly sourced from the paper, formatted according to the task type (MCQ, numerical,
835 or free-form). Experts additionally curated background knowledge necessary for informed prediction,
836 selecting relevant domain principles, previously established findings, and theoretical frameworks
837 from the source papers or from expert domain knowledge. Fig. 2 provides a representative example of
838 this extraction and background curation process.

839 B.4 HUMAN BASELINE RECRUITMENT DETAILS

840
841 In addition to the experts involved in benchmark construction, we recruited a separate group of experts
842 to serve as human baseline subjects. These participants were selected to represent an expert-level
843 baseline for the prediction tasks. Human baseline subjects were presented with benchmark questions
844 and asked to provide an answer, explain their reasoning, and report their confidence. To mirror the
845 evaluation protocol used for LLM baselines, each subject completed a second round of the same
846 questions after being provided with the curated background knowledge associated with the task.

847 The human baseline cohort consists primarily of domain experts, with 74.4% holding doctoral degrees,
848 17.9% holding master’s degrees, and 7.7% holding bachelor’s degrees. In terms of primary area of
849 expertise, 48.7% specialize in biology, 33.3% in chemistry, and 17.9% in physics. The cohort also
850 reflects broad geographic diversity, including participants from the United States (33.3%), Argentina
851 (17.9%), the United Kingdom (15.4%), Mexico (7.7%), and Colombia (5.1%). Fig. 12 provides a
852 detailed demographic breakdown.

853 To ensure that human baseline performance reflects expert-level reasoning rather than domain
854 mismatch, we performed a rigorous assignment process aligning each subject’s area of expertise with
855 the corresponding task subdomains. The resulting expertise-to-task mapping is summarized in Tab. 1.
856

857 B.5 QUALITY CONTROL DETAILS

858
859 All data undergoes a multi-stage review process to ensure scientific rigor. Initial screening filters
860 questions where the first version of the paper appeared online on or before March 31, 2025, experiments
861 are simulations or theoretical derivations, answers are directly stated in experimental setup descriptions,
862 phrasing is ambiguous, required predictions exceed available information, or ground truth conflicts
863 with source papers. Questions passing initial screening go through two layers of domain expert
reviewers who verify experimental setup precision sufficiency for informed reasoning, background

864 knowledge necessity and sufficiency, ground truth clarity and proper sourcing, and appropriate
865 difficulty level.

866 For multiple-choice questions, reviewers ensure distractors represent plausible alternatives arising
867 from reasonable but incorrect assumptions rather than obviously wrong options. For free-form
868 questions, reviewers confirm that evaluation rubrics capture essential scientific reasoning without
869 being overly prescriptive about phrasing, and that rubric criteria are mutually exclusive and collectively
870 exhaustive, with each criterion validated to a binary outcome. For numerical value questions, reviewers
871 verify acceptable ranges are neither unrealistically narrow nor trivially broad, reflecting realistic
872 experimental measurement precision and variability. Questions flagged during review undergo
873 revision or removal if fundamental problems cannot be resolved.

875 B.6 DATA DIVERSITY DETAILS

877 The benchmark spans 33 specialized subdomains across physics, biology, and chemistry, ensuring
878 models encounter the full spectrum of experimental reasoning required in modern scientific practice.
879 Within physics, questions draw from 9 subdomains such as experimental condensed matter physics,
880 quantum and atomic physics, and high energy particle physics. Biology questions cover 14 subdomains
881 such as molecular biology, neuroscience, plant biology, and ecology. Chemistry spans 10 subdomains
882 such as organic chemistry, catalysis, and polymer chemistry.

883 Question complexity varies systematically along multiple axes. Experimental systems range from
884 controlled laboratory setups with few interacting components to complex biological systems with
885 emergent properties. Some questions require single-step causal reasoning, while others demand
886 multi-hop inference chains such as integrating thermodynamics, kinetics, and material properties.
887 Background knowledge requirements span a continuum from questions answerable via undergraduate-
888 level principles to those requiring specialized domain expertise typically held only by active researchers
889 in the relevant subdomain.

890 Domain distribution remains balanced to prevent overfitting to particular experimental contexts,
891 with 25% of questions from physics, 50% from biology, and 25% from chemistry. Question format
892 distribution is similarly controlled, with 40% multiple-choice, 32% free-form, and 28% numerical
893 value questions. Together, these diversity dimensions ensure the benchmark probes models' general
894 capacity for experimental outcome prediction rather than narrow pattern matching on specific
895 experimental templates or domain conventions.

896 B.7 HUMAN BASELINE EXPERT - TASK SUBDOMAIN MAPPING

898 Table 1: Subfield expertise of human annotators, grouped by the task domains (Physics, Chemistry,
899 Biology) and subdomains.

901 Task Domain	902 Subdomain	903 Human Baseline Subfields
904 Physics	905 All Physics	906 Advanced Chemical Engineering, Applied And Interdisciplinary Physics, Applied Physics And Interdisciplinary, Chemical Engineering, Classical And Mechanical Physics, Condensed Matter And Materials, Electromagnetism And Optics, Engineering Physics, High-energy And Nuclear Physics, Radiophysics & Electronics, Theoretical Physics, Zoology
	907 Condensed Matter & Materials Physics	908 Advanced Chemical Engineering, Applied Physics And Interdisciplinary, Chemical Engineering, Condensed Matter And Materials, Electromagnetism And Optics, Engineering Physics, Radiophysics & Electronics
	909 Materials Chemistry	910 Condensed Matter And Materials, Engineering Physics
	911 Optics, Photonics & Laser Physics	912 Applied Physics And Interdisciplinary, Condensed Matter And Materials, Electromagnetism And Optics, Engineering Physics, Radiophysics & Electronics, Zoology
	913 High-Energy / Nuclear / Particle Physics	914 Engineering Physics, High-energy And Nuclear Physics, Radiophysics & Electronics, Theoretical Physics, Zoology
	915 Applied & Instrumentation Physics	916 Applied And Interdisciplinary Physics, Applied Physics And Interdisciplinary, Classical And Mechanical Physics, Condensed Matter And Materials, Electromagnetism And Optics, Engineering Physics, High-energy And Nuclear Physics, Radiophysics & Electronics
	917 Quantum & Atomic Physics	918 Applied Physics And Interdisciplinary, Condensed Matter And Materials, Electromagnetism And Optics, Engineering Physics, Radiophysics & Electronics, Zoology
	919 Plasma & Nonlinear Physics	920 Applied Physics And Interdisciplinary, Classical And Mechanical Physics, Electromagnetism And Optics, Engineering Physics, Radiophysics & Electronics

Task Domain	Subdomain	Human Baseline Subfields		
918 919 920 921 922 923	Biophysics	Advanced Chemical Engineering, Applied Physics And Interdisciplinary, Chemical Engineering, Condensed Matter And Materials, Electromagnetism And Optics, Radiophysics & Electronics		
	Mechanical / Energy / Thermo / Fluid Physics	Classical And Mechanical Physics, Condensed Matter And Materials, Engineering Physics, Radiophysics & Electronics		
924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950	Chemistry	All Chemistry		
		Analytical Chemistry	Advanced Chemical Engineering, Analytical Chemistry, Antimicrobial Resistance, Bio-organic Chemistry, Biochemistry And Molecular Biology, Chemical Biology, Chemical Engineering, Chemical Sciences, Digital Technologies Applied To Education, Electrochemistry, Engineering Physics, Materials And Inorganic Chemistry, Molecular And Cellular Biology, Molecular Biology And Genetics, Organic And Biological Chemistry, Principles Of Biochemistry, Pure Chemistry, Zoology	
		Materials Chemistry	Analytical Chemistry, Bio-organic Chemistry, Biochemistry And Molecular Biology, Chemical Biology, Chemical Engineering, Digital Technologies Applied To Education, Electrochemistry, Materials And Inorganic Chemistry, Organic And Biological Chemistry	
		Catalysis	Biochemistry, Biochemistry And Molecular Biology, Catalysis And Environmental Chemistry, Chemical Biology, Chemical Engineering, Chemical Sciences, Digital Technologies Applied To Education, Electrochemistry, Green Chemistry, Materials And Inorganic Chemistry, Principles Of Biochemistry, Pure Chemistry	
		Physical Chemistry	Advanced Chemical Engineering, Analytical Chemistry, Chemical Engineering, Chemical Sciences, Digital Technologies Applied To Education, Materials And Inorganic Chemistry, Organic And Biological Chemistry, Principles Of Biochemistry, Pure Chemistry	
		Organic Chemistry	Analytical Chemistry, Bio-organic Chemistry, Biochemistry And Molecular Biology, Catalysis And Environmental Chemistry, Chemical Biology, Chemical Engineering, Digital Technologies Applied To Education, Electrochemistry, Materials And Inorganic Chemistry, Organic And Biological Chemistry, Zoology	
		Nanotechnology / Nanochemistry	Analytical Chemistry, Biochemistry, Biochemistry And Molecular Biology, Catalysis And Environmental Chemistry, Chemical Biology, Chemical Engineering, Digital Technologies Applied To Education, Electrochemistry, Green Chemistry, Materials And Inorganic Chemistry, Organic And Biological Chemistry, Principles Of Biochemistry, Pure Chemistry	
		Biochemistry	Antimicrobial Resistance, Biochemistry, Electrochemistry, Molecular And Cellular Biology, Molecular Biology And Genetics, Organic And Biological Chemistry, Principles Of Biochemistry, Pure Chemistry	
		Inorganic Chemistry	Analytical Chemistry, Catalysis And Environmental Chemistry, Materials And Inorganic Chemistry	
		Environmental Chemistry	Advanced Chemical Engineering, Analytical Chemistry, Chemical Engineering, Materials And Inorganic Chemistry, Zoology	
		Polymer Chemistry	Chemical Engineering, Digital Technologies Applied To Education, Materials And Inorganic Chemistry, Organic And Biological Chemistry	
	951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971	Biology	All Biology	
			Microbiology	Antimicrobial Resistance, Biochemistry, Biological Engineering, Biological Sciences, Biomedical Engineering, Biomedical Sciences, Cell Biology, Chemical Engineering, Ecology, Microbiology, Microbiology And Cell Science, Molecular And Cellular Biology, Molecular Biology And Genetics, Neurobiology And Behavior, Software Engineering, Systems And Synthetic Biology, Taxonomy And Biodiversity
			Cancer Biology / Oncology	Antimicrobial Resistance, Biochemistry, Biological Engineering, Biological Sciences, Biomedical Engineering, Biomedical Sciences, Cell Biology, Chemical Engineering, Clinical Drug Development, Genetics, Immunology, Microbiology And Cell Science, Molecular And Cellular Biology, Molecular Biology, Molecular Biology And Genetics, Research And Data Analysis, Software Engineering, Taxonomy And Biodiversity
		Neuroscience / Neurobiology	Antimicrobial Resistance, Biochemistry, Biological Engineering, Biomedical Engineering, Cell Biology, Chemical Engineering, Clinical Drug Development, Developmental Biology, Genetics, Immunology, Molecular And Cellular Biology, Molecular Biology, Molecular Biology And Genetics, Neurobiology And Behavior, Physiology, Systems And Synthetic Biology	
		Ecology	Biochemistry, Biological Engineering, Biological Sciences, Biomedical Engineering, Biomedical Sciences, Cell Biology, Chemical Engineering, Ecology, Genetics, Microbiology, Microbiology And Cell Science, Observational Oceanography, Plant Sciences, Research And Data Analysis, Systems And Synthetic Biology, Taxonomy And Biodiversity	
		Immunology	Bio-organic Chemistry, Biochemistry, Biological Engineering, Biomedical Engineering, Biomedical Sciences, Chemical Engineering, Immunology, Microbiology And Cell Science, Software Engineering, Systems And Synthetic Biology, Zoology	

Task Domain	Subdomain	Human Baseline Subfields
972	Molecular Biology	Antimicrobial Resistance, Bio-organic Chemistry, Biochemistry, Biological Engineering, Biological Sciences, Biomedical Engineering, Biomedical Sciences, Cell Biology, Chemical Engineering, Genetics, Microbiology And Cell Science, Molecular And Cellular Biology, Molecular Biology, Molecular Biology And Genetics, Research And Data Analysis, Software Engineering, Taxonomy And Biodiversity
973		
974	Pharmacology / Toxicology	Biochemistry, Biological Sciences, Biomedical Sciences, Cell Biology, Clinical Drug Development, Genetics, Immunology, Microbiology And Cell Science, Observational Oceanography, Physiology, Research And Data Analysis, Software Engineering
975		
976	Plant Biology	Biochemistry, Biological Sciences, Developmental Biology, Ecology, Genetics, Observational Oceanography, Plant Sciences, Research And Data Analysis, Systems And Synthetic Biology, Taxonomy And Biodiversity
977		
978	Animal Behavior	Biochemistry, Biological Sciences, Cell Biology, Clinical Drug Development, Developmental Biology, Genetics, Microbiology, Molecular Biology, Observational Oceanography, Physiology, Systems And Synthetic Biology, Taxonomy And Biodiversity, Zoology
979		
980	Cell Biology	Antimicrobial Resistance, Bio-organic Chemistry, Biochemistry, Biological Engineering, Biological Sciences, Biomedical Engineering, Biomedical Sciences, Cell Biology, Chemical Engineering, Clinical Drug Development, Developmental Biology, Genetics, Immunology, Microbiology And Cell Science, Molecular And Cellular Biology, Molecular Biology, Molecular Biology And Genetics, Neurobiology And Behavior, Physiology, Research And Data Analysis, Software Engineering, Taxonomy And Biodiversity
981		
982	Physiology	Biochemistry, Biological Engineering, Biological Sciences, Biomedical Engineering, Biotechnology, Cell Biology, Chemical Engineering, Clinical Drug Development, Genetics, Microbiology, Molecular And Cellular Biology, Molecular Biology, Neurobiology And Behavior, Observational Oceanography, Physiology, Plant Sciences, Systems And Synthetic Biology, Taxonomy And Biodiversity
983		
984	Biochemistry	Biochemistry, Biochemistry And Molecular Biology, Biological Engineering, Biomedical Engineering, Cell Biology, Chemical Biology, Chemical Engineering, Clinical Drug Development, Genetics, Molecular Biology, Physiology, Software Engineering, Zoology
985		
986	Genetics	Biochemistry, Biological Sciences, Biomedical Sciences, Cell Biology, Clinical Drug Development, Genetics, Microbiology, Microbiology And Cell Science, Molecular Biology, Observational Oceanography, Plant Sciences, Systems And Synthetic Biology, Taxonomy And Biodiversity
987		
988	Bioengineering / Biomaterials	Antimicrobial Resistance, Biochemistry, Biological Sciences, Biomedical Sciences, Cell Biology, Green Chemistry, Microbiology And Cell Science, Molecular And Cellular Biology, Molecular Biology, Molecular Biology And Genetics, Observational Oceanography, Physiology, Systems And Synthetic Biology
989		
990		
991		
992		
993		
994		
995		
996		
997		
998		
999		
1000		
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010		
1011		
1012		
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020		
1021		
1022		
1023		
1024		
1025		

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Table 2: Task distribution by scientific subfield: number of tasks per Biology, Physics, and Chemistry subdomain.

Field	Subfield	Count
Physics	Condensed Matter & Materials Physics	33
	Materials Chemistry	17
	Optics, Photonics & Laser Physics	16
	High-Energy / Nuclear / Particle Physics	15
	Applied & Instrumentation Physics	13
	Quantum & Atomic Physics	10
	Plasma & Nonlinear Physics	5
	Biophysics	3
	Mechanical / Energy / Thermo / Fluid Physics	2
Chemistry	Analytical Chemistry	18
	Materials Chemistry	17
	Catalysis	16
	Physical Chemistry	14
	Organic Chemistry	13
	Nanotechnology / Nanochemistry	10
	Biochemistry	8
	Inorganic Chemistry	6
	Environmental Chemistry	4
	Polymer Chemistry	3
Biology	Microbiology	36
	Cancer Biology / Oncology	28
	Neuroscience / Neurobiology	19
	Ecology	17
	Immunology	16
	Molecular Biology	14
	Pharmacology / Toxicology	13
	Plant Biology	13
	Animal Behavior	13
	Cell Biology	10
	Physiology	9
	Biochemistry	8
	Genetics	8
Bioengineering / Biomaterials	3	

1080 C EXAMPLE DATA

1081

1082

1083 C.1 TASK EXAMPLES

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

Physics: Free-Form Question

Paper Title: Compact Continuous Cold Atomic Beam from a Single Cell with 3D Cooling and Ultra-low Light Shift

Link to The Paper: <https://arxiv.org/abs/2510.13126>

Experimental Setup: Researchers investigated a compact single-cell source of a continuous cold-atom beam (^{87}Rb) that achieves simultaneous 3D cooling by integrating a two-dimensional magneto-optical trap (2D MOT) with an off-axis moving optical molasses (OM). A vapor-cell apparatus (overall length ≈ 170 mm) provided transverse MOT cooling with circularly polarized beams detuned by $\Delta\text{MOT} = -4\Gamma$ from the $F = 2 \rightarrow F' = 3$ D_2 transition and a cylindrical quadrupole field (≈ 10 G cm^{-1}), where Γ is the natural linewidth. Longitudinal cooling and velocity control were realized with two pairs of lin \perp lin OM beams oriented 20° to the extraction axis, detuned by $\Delta\text{OM} = -5\Gamma$ and symmetrically shifted by $\pm\delta\text{OM}$ to set the mean atomic speed ($\approx 5\text{--}20$ m s^{-1}) over an OM interaction length $l_{\text{OM}} \approx 50$ mm. Custom in-vacuum mirrors formed the off-axis geometry and incorporated a 0.8 mm output aperture to collimate the beam (cooling length $l_c \approx 50$ mm) while suppressing near-resonant stray light. The setup included permanent-magnet field generation, state-preparation “plug” lasers 40 mm downstream for sharp time-of-flight (TOF) edges, and fluorescence detection at 294 mm with a calibrated photomultiplier tube (PMT) to extract longitudinal temperature, velocity, and flux. For coherence diagnostics, two $\pi/2$ Raman beams separated by $L = 100$ mm in a magnetically shielded region produced spatial-domain Raman–Ramsey fringes, enabling quantification of decoherence and ultra-low light shift (typ. -0.51 Hz) under operating MOT power.

Measurements Taken:

- Time-of-flight (TOF) time series and distribution obtained from the emitted fluorescence from the atoms in $F=2$ state, collected with imaging optics and recorded by a calibrated PMT at a primary detection distance of 294 mm.

Outcome Prediction Question: Researchers investigated the longitudinal temperature and atomic flux of a continuous cold ^{87}Rb beam using a time-of-flight (TOF) method. The temperature was extracted from the FWHM of the TOF distribution, while the flux was obtained from the integrated spectral density. Based on measurements for a saturation intensity of 1.67 mW/cm 2 , what outcome would researchers expect for the change in longitudinal temperature and atomic flux when the MOT power is increased?

Ground Truth Answer: Increasing MOT power raises the flux but affects the temperature only weakly.

Background Knowledge:

- Combining a 2D MOT with an off-axis moving OM yields a high-flux beam with significantly reduced longitudinal temperature compared to conventional MOT-based sources.
- Continuous operation of cold-atom beam sources eliminates the dead time inherent to pulsed sources and thus suppresses aliasing noise from undersampling.

Rubrics:

- Response states that increasing the magneto-optical trap power increases atomic flux.
- Response states that increasing the magneto-optical trap power has a little influence on temperature.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Physics: Multiple-Choice Question

Paper Title: Ionization and temperature measurements in warm dense copper using x-ray absorption spectroscopy

Link to The Paper: <https://arxiv.org/abs/2509.13272>

Experimental Setup: Researchers investigated the ionization and temperature of warm dense copper (Cu) using X-ray absorption spectroscopy (XAS) at the OMEGA Laser Facility to characterize plasmas at several times solid density. The experimental configuration consists of a planar target and a separate backlighter positioned 3 mm away. A series of 60 laser beams, delivering 3.4–5.4 kJ per side of 351 nm light, and the achieved laser intensity is 161 - 770 TW/cm² over the three pulse length configurations, was symmetrically focused onto a planar buried-layer target composed of 125 μm CH ablators enclosing a 10 μm-thick Cu foil (8.96 g/cm³ solid density) with a 500 μm diameter, surrounded by an Au washer. The laser spot (≈ 880 μm diameter) was smoothed with distributed phase plates and spectral dispersion to generate uniform counter-propagating shocks. A 6 μm Ge backlighter foil, coated on graphite and irradiated with six additional beams (≈1.2 kJ, 500 ps pulse), is produced at a spot diameter of 140 μm. The transmitted x-rays were recorded using the EFX flat-crystal spectrometer (Si 111) over the 6.3–11.4 keV range on an image plate with Mn, Fe, and W filters serving as fiducial markers. Shock timing and planarity, as well as shock break-in and break-out of the Cu layer, were verified through a line-imaging VISAR system and a streaked optical pyrometer (SOP) on one-sided targets, ensuring symmetric compression and precise backlighter synchronization. 3 VISAR measurement is done with 1 ns, 2 ns, or 3 ns square pulses using 14 beams per side, respectively. Each measurement has two VISAR channels with different sensitivities; one leg was set with 33.66 μm/ns/fringe, and the second with 13.538 μm/ns/fringe.

Measurements Taken:

- Shock breakout times (in ns) and planarity were measured with the VISAR system.
- Shock velocity time history as a function of position across the target measured with the VISAR system.

Outcome Prediction Question: An investigation into shock breakout times and shock velocity time histories as a function of position across the target of warm dense copper (Cu) plasma is conducted using a VISAR system. The experimental configuration consists of a planar target and a separate backlighter positioned 3 mm away. A series of 60 laser beams was symmetrically focused onto a planar buried-layer target surrounded by an Au washer. The laser spot was smoothed with distributed phase plates and spectral dispersion to generate uniform counter-propagating shocks, compressing the Cu layer. A Ge backlighter foil, coated on graphite and irradiated with six additional beams, is produced. The transmitted X-rays were also recorded using the EFX flat-crystal spectrometer. Which behavior is most likely observed?

- A. Shocks were non-planar over the target region, and warm dense copper shows Ionization Potential Depression (IPD).
- B. Shocks were highly planar over the target region, and the absorption spectra of warm dense copper features blue shift of both the K-edge and the bound-bound resonance 1s→3p absorption relative to the cold edge.
- C. Shocks were highly planar over the target region, and the absorption spectra of warm dense copper features red shift of both the K-edge and the bound-bound resonance 1s→3p absorption relative to the cold edge.
- D. Shocks were highly planar over the target region, and the absorption spectra of warm dense copper features blue shift of the K-edge relative to the cold edge, but no shift for the bound-bound resonance 1s→3p absorption.

Ground Truth Answer: B

Background Knowledge:

- Generating warm dense matter in the laboratory often involves significant temporal and spatial gradients that complicate the analysis of experimental observables. Incorporating gradients in the analysis of experimental data, while possible, increases the uncertainties in the inferred plasma conditions.
- At these high-density conditions, the measured Cu K-edge exhibits sensitivity to the electron temperature, allowing for a direct inference of the temperature from the slope of the Cu K-edge.
- Temperature sensitivity of the K-edge can still be the dominant edge effect, in general, as the temperature nears the Fermi energy, the K-edge shape of the non-degenerate material becomes unsuitable as a temperature inference.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Physics: Numerical Value Question

Paper Title: A sub-volt near-IR lithium tantalate electro-optic modulator

Link to The Paper: <https://arxiv.org/abs/2505.00906>

Experimental Setup: Researchers fabricated a TFLT MZM operating at a near-IR wavelength of 737 nm. The fabricated unbalanced MZM consists of a directional coupler as an input beamsplitter and a $L = 5$ mm long electrode in the ground-signal-ground configuration, followed by another directional coupler at the output. Grating couplers are used to couple light on and off the chip to near-IR single-mode fibers. The optical layer of the device is defined using 150 keV electron-beam lithography with 500 nm-thick ma-N2405 resist on top of a 200 nm-thick x-cut TFLT-on-SiO₂ layer. The waveguide width is designed to be 600 nm. The SiO₂ layer is 2 μm -thick and is on a Si substrate. The TFLT is etched by 100 nm using an Ar⁺-based inductively coupled plasma reactive ion etching. Etch-induced re-deposition is removed using a high-pH solution. The devices are then annealed in an O₂ atmosphere at 520°C for 2 h to mitigate etch-induced imperfections. For the MZMs, an 800 nm-thick SiO₂ cladding layer is then deposited by plasma-enhanced chemical vapor deposition. The DC bias stability of two electro-optic Mach-Zehnder modulators is compared. The first modulator is fabricated using thin-film lithium tantalate (TFLT), and the second, serving as a counterpart, is fabricated with a similar process using thin-film lithium niobate (TFLN). For the test, each modulator is subjected to a constant on-chip optical power of 4.3 dBm at a wavelength of 737 nm. A DC step voltage is applied to each device to set its operating point at quadrature bias. The output optical power from the modulator is then monitored over 16 minutes in ambient conditions to measure any drift from this bias point. To measure the DC bias stability of MZM over long timescales. First, it applied a 0.1 Hz-frequency square wave to the modulator using an on-chip optical power, and measured the modulator response with a photodetector.

Measurements Taken:

- The output optical power as a function of time over 16 minutes for the TFLT modulator.
- The output optical power as a function of time over 16 minutes for the TFLN modulator.
- The total DC bias drift, in decibels (dB), for the TFLT modulator.
- The total DC bias drift, in decibels (dB), for the TFLN modulator.

Outcome Prediction Question: An experiment compares the long-term stability of two Mach-Zehnder modulators, one made from thin-film lithium tantalate (TFLT) and a counterpart from thin-film lithium niobate (TFLN). Both are operated with 4.3 dBm of on-chip optical power at 737 nm and biased at quadrature. The output power is monitored for 16 minutes to quantify the DC bias drift. To measure the DC bias stability of MZM over long timescales. First, it applied a 0.1 Hz-frequency square wave to the modulator using an on-chip optical power, and measured the modulator response with a photodetector. Based on the experimental results, what is the total measured DC bias drift, in decibels (dB), for the thin-film lithium niobate (TFLN) modulator?

Ground Truth Answer: $\Delta\text{DC bias drift} = [7.2\text{--}8.8]$ dB at 16 min for the TFLN modulator operated at 4.3 dBm optical power (737 nm). No CI/SE/SD reported \rightarrow fallback ± 0.8 dB applied.

Background Knowledge:

- In particular, the relaxation rate will increase with more applied optical power and can be exacerbated with applied DC or RF field. This effect reduces the DC stability of electro-optic circuits, such as Mach-Zehnder modulators (MZMs), and has been one of the main challenges faced by TFLN photonics

Biology: Free-Form Question

Paper Title: Dopamine induces fear extinction by activating the reward-responding amygdala neurons

Link to The Paper: <https://pmc.ncbi.nlm.nih.gov/articles/PMC12067255/>

Experimental Setup: Researchers tested whether ventral tegmental area (VTA) dopamine signaling in the basolateral amygdala (BLA) drives fear extinction by acting on reward-responding posterior BLA (pBLA) neurons versus fear-coding anterior BLA (aBLA) neurons, using adult mice (DAT-IRES-Cre; EYFP controls; subtype mapping with Rspo2-Cre for aBLA and Ppp1r1b/Cartpt-Cre for pBLA). DAT-Cre mice received bilateral VTA injections of Cre-dependent ChR2-EYFP (activation) or eNpHR3.0-EYFP (inhibition); controls received EYFP; optic fibers were implanted over pBLA or aBLA to manipulate VTA→BLA terminals. Training: Day 1 contextual fear conditioning (baseline ~3 min, then 3 footshocks, 0.60 mA, 2 s); Day 2 45-min extinction (no shocks); Day 3 10-min retrieval. Intervention (extinction only): starting 5 min into extinction, deliver 8 cycles of 3-min light separated by 2-min no-light (activation: blue 450–470 nm, 8–12 mW, 20 Hz pulses; inhibition: green 520–550 nm, 8–12 mW, continuous) with fibers targeted to pBLA or aBLA. Behavior videos were recorded with VideoFreeze software and freezing level was scored manually by experimenters who were blinded to conditions or automatically with DeepLabCut behavior analysis toolbox and custom Python code (68). Freezing was quantified in 5-min bins across extinction and again during retrieval.

Measurements Taken:

- Extinction learning: Percent freezing per 5-min bin across the 45-min Day 2 session (9 bins). Scored manually by experimenters who were blinded to conditions or automatically with DeepLabCut behavior analysis toolbox and custom Python code (68).
- Extinction memory: Percent freezing during the Day 3 retrieval test (10 min). Scored manually by experimenters who were blinded to conditions or automatically with DeepLabCut behavior analysis toolbox and custom Python code (68).

Outcome Prediction Question: Mice underwent contextual fear conditioning (Day 1: context + three 0.60 mA, 2 s shocks), 45-min extinction (Day 2, no shocks), and 10-min retrieval (Day 3). During extinction, VTA dopamine terminals in pBLA (Ppp1r1b⁺) or aBLA (Rspo2⁺) were optogenetically manipulated beginning 5 min into the session using 8 cycles of 3 min light separated by 2 min: activation (blue 450–470 nm, 8–12 mW, 20 Hz) or inhibition (green 520–550 nm, 8–12 mW, constant). Freezing was binned in 5-min windows across extinction and measured again at retrieval. How do these projection-specific manipulations (activation and inhibition of VTA dopamine terminals in the pBLA and in aBLA) affect fear extinction and retrieval compared with EYFP controls?

Ground Truth Answer: Activation of VTA dopamine terminals in the pBLA promotes faster extinction and improved retrieval, indicating an enhancement of extinction learning. In contrast, inhibition of pBLA dopamine input impairs both extinction and retrieval. Activation of VTA terminals in the aBLA leads to increased freezing later in extinction and poorer retrieval performance, suggesting interference with extinction memory formation, while inhibition of aBLA terminals produces no reliable behavioral change.

Background Knowledge:

- Fear extinction is a form of new learning that allows for the adaptive control of fear behaviors and is commonly studied using Pavlovian conditioning tasks.
- aBLA Rspo⁺ neurons encode negative valence and drive aversive behaviors whereas pBLA Ppp1r1b⁺ neurons encode positive valence and drive appetitive behaviors.
- VTA dopamine as a teaching signal: DA activity to shock omission can initiate extinction learning and is required for extinction.
- Terminal activation (ChR2, blue, pulsed) vs inhibition (eNpHR3.0, green, constant) at BLA terminals tests sufficiency/necessity of VTA→BLA pathways.
- Freezing is the behavioral measure; decreases across 5-minute bins and at retrieval indicate successful extinction.

Rubrics:

- The response should state that activation of ventral tegmental area dopamine terminals in the posterior basolateral amygdala of adult mice promotes faster extinction compared to control. Use of acronyms such as VTA or pBLA are acceptable.
- The response should state that activation of ventral tegmental area dopamine terminals in the posterior basolateral amygdala of adult mice improves retrieval compared to control. Use of acronyms such as VTA or pBLA are acceptable.
- The response should state that inhibition of ventral tegmental area dopamine terminals in the posterior basolateral amygdala of adult mice impairs extinction compared to control. Use of acronyms such as VTA or pBLA are acceptable.
- The response should state that inhibition of ventral tegmental area dopamine terminals in the posterior basolateral amygdala of adult mice impairs retrieval compared to control. Use of acronyms such as VTA or pBLA are acceptable.
- The response should state that activation of ventral tegmental area terminals in the anterior basolateral amygdala of adult mice leads to increased freezing later in extinction compared to control. Use of acronyms such as VTA or aBLA are acceptable.
- The response should state that activation of ventral tegmental area terminals in the anterior basolateral amygdala of adult mice leads to poorer retrieval performance compared to control. Use of acronyms such as VTA or aBLA are acceptable.
- The response should state that inhibition of ventral tegmental area terminals in the anterior basolateral amygdala of adult mice produces no reliable behavioral change compared to control. Use of acronyms such as VTA or aBLA are acceptable.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Biology: Multiple-Choice Question

Paper Title: Social Tolerance and Innovation in Capuchins: socially more tolerant brown capuchins are better problem-solvers than less tolerant white-faced capuchins

Link to The Paper: <https://www.biorxiv.org/content/10.1101/2025.09.05.674457v1.full>

Experimental Setup: Researchers tested three groups of white-faced capuchins (*Cebus capucinus*) (n = 23 individuals in total) and three groups of brown capuchins (*Sapajus apella*) (n = 20 individuals in total) to explore and compare the relationship between social tolerance and problem-solving propensities. To measure social tolerance, they prepared an area of 1 m² per five animals in the group, in which they distributed apple pieces and measured the proportion of individuals within the co-feeding area at each scan sample. To measure problem-solving propensities, they designed three versions of novel extractive foraging devices requiring one to three steps to acquire the food reward. For the first puzzle, animals had to rotate a door to either the left or right to access a hidden reward (1/24 of an apple) by reaching into a box. For the second puzzle, animals had to pull on a chain reaching out of a box, which moved a blockade out of the way so that they could push in a door and reach into the box. For the third puzzle, animals had to pull a metal rod blocking a slider that had to be pulled upwards and held in position to reach into the box and then pull on a chain to access the hidden reward. Researchers analyzed the approaching, exploring, and solving behaviour separately.

Measurements Taken:

- Proportion of individuals within the co-feeding area at each scan sample (social tolerance)
- Proportion of individuals within the puzzle area at each scan sample (social tolerance)
- Number of approaches to a food puzzle area
- Approaching a food puzzle area duration
- Approaches to a food puzzle area latency
- Number of exploration events (touch, sniff, interact) during the approaches to a food puzzle area
- Number of times the capuchins successfully solved the puzzles
- Exploration of food puzzle events latency
- Time to solve a puzzle

Outcome Prediction Question: Researchers tested three groups of white-faced capuchins (*Cebus capucinus*) and three groups of brown capuchins (*Sapajus apella*) to explore and compare the relationship between social tolerance and problem-solving propensities. To measure social tolerance, they prepared an area of 1 m² per five animals in the group, in which they distributed apple pieces and measured the proportion of individuals within the co-feeding area at each scan sample. To measure problem-solving propensities, they designed three versions of novel extractive foraging devices requiring one to three steps to acquire the food reward. Which of the following outcomes is most likely?

- A. Both species should show the same levels of social tolerance and problem-solving propensities.
- B. White-faced capuchins should show the highest level of social tolerance and problem-solving propensities.
- C. White-faced capuchins should show the lowest level of social tolerance and problem-solving propensities.
- D. White-faced capuchins should show the highest level of social tolerance and the lowest level of problem-solving propensities.

Ground Truth Answer: C

Background Knowledge:

- Social tolerance has increasingly been linked to the facilitation of social learning across a variety of species, including chimpanzees, orangutans, macaques, capuchin monkeys, lemurs, and birds.
- White-faced capuchins (*Cebus capucinus*) and brown capuchins (*Sapajus apella*) exhibit a diverse array of traditions.
- White-faced capuchins (*Cebus capucinus*) are less known for using tools (but see Barrett et al., 2018), but they regularly engage in object use (Boinski, 1988).
- Robust capuchins (*Sapajus* spp.) have fewer documented social traditions but exhibit a wide range of foraging traditions, including tool-use, and show notable social tolerance in these contexts, tolerating close proximity of conspecifics.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Biology: Numerical Value Question

Paper Title: GsMTx4-loaded GelMA promotes tendon regeneration and suppresses heterotopic ossification via the Apelin signaling pathway

Link to The Paper:

<https://www.sciencedirect.com/science/article/pii/S0142961225004260?via%3Dihub>

Experimental Setup: Researchers employed Male Sprague Dawley (SD) rats (10–12 weeks old, weighing 250–300 g) as animal model for studying tendon repair and regeneration. A central defect (1 mm in width and 5 mm in length) was created in the Achilles tendon using two parallel No.15 surgical blades. Subsequently, the skin was sutured using 4-0 Vicryl sutures. The rats received temgesic (0.3 mg/kg of body weight) for three consecutive days following the surgery to manage pain. The rats were randomly assigned to one of four groups: Achilles tendon defect (ATD) (no treatment), GelMA, GelMA + 50 μg GsMTx4, GelMA + 100 μg GsMTx4. At the time of injury, a mixture of GelMA and LAP (Lithium Phenyl-2,4,6-Trimethylbenzoylphosphinate) (20 μl), loaded with 50 or 100 μg GsMTx4 where appropriate, was placed within the ATD of treated animals and transformed into the gel state with a blue light source (3 W, 405 nm) for 30 s at a distance of 2 cm from the defect. These animals were euthanized at 2, 4, and 8 weeks post-treatment, with six rats per group per time point. The harvested Achilles tendons were fixed in 4% paraformaldehyde at room temperature for 24 h. Following fixation, the samples were rinsed with running water and dehydrated with an ethanol gradient, and embedded in paraffin. The blocks were sectioned at 5 μm thickness using a microtome and stained with Hematoxylin and Eosin (H&E). Semi-quantitative analysis of H&E staining results was conducted according to the modified Bonar score.

Measurements Taken:

- Histologic Bonar Score (ATD, GelMA, GelMA + 50 μg GsMTx4, GelMA + 100 μg GsMTx4): 2 weeks; 4 weeks; 8 weeks.

Outcome Prediction Question: Researchers employed Male Sprague Dawley (SD) rats (10–12 weeks old, weighing 250–300 g) as animal model for studying tendon repair and regeneration. A central defect (1 mm in width and 5 mm in length) was created in the Achilles tendon using two parallel No.15 surgical blades. The rats were randomly assigned to one of four groups: Achilles tendon defect (ATD) (no treatment), GelMA, GelMA + 50 μg GsMTx4, GelMA + 100 μg GsMTx4. At the time of injury, a mixture of GelMA and LAP (Lithium Phenyl-2,4,6-Trimethylbenzoylphosphinate) (20 μl), loaded with 50 or 100 μg GsMTx4 where appropriate, was placed within the ATD of treated animals and transformed into the gel state with a blue light source. The animals were euthanized at 2, 4, and 8 weeks post-treatment. The harvested Achilles tendons were embedded in paraffin, sectioned using a microtome, and stained with Hematoxylin and Eosin (H&E). Semi-quantitative analysis of H&E staining results was conducted according to the modified Bonar score (BS). Based on the reported values of the BS for Achilles tendon repair and regeneration, what is the predicted difference of the BS (in points) between the GelMA and the GelMA + 100 μg GsMTx4 groups 8-weeks post treatment?

Ground Truth Answer: Δ BS (GelMA - GelMA + 100 μg GsMTx4) 8-weeks post treatment = 4 - 6 points; derived from BS GelMA 8-weeks post treatment = ~9 points, BS GelMA + 100 μg GsMTx4 8-weeks post treatment = ~4 points. Note: No CI/SE/SD reported -> fallback \pm 10% units (rounded) applied.

Background Knowledge:

- Tendon regeneration is highly relied on the surrounding mechanical environment.
- Studies have demonstrated the importance of Piezo1 in modulating cellular behaviors to mechanical cues, such as cell migration, differentiation, proliferation, and extracellular matrix synthesis.
- GelMA hydrogel demonstrates excellent biocompatibility and sustained release properties.
- The mechanosensitive ion channel Piezo1 is inhibited by the peptide GsMTx4

Chemistry: Free-Form Question

Paper Title: An investigation of the physical and chemical changes of Pd nanoparticles on carbon supports in response to the release of hydrogen from aqueous formate solutions

Link to The Paper:

<https://chemrxiv.org/engage/chemrxiv/article-details/68d16d29f2aff16770fa93bd>

Experimental Setup: Researchers prepared and analyzed Pd nanoparticles supported on carbon materials to examine their structural and chemical evolution during hydrogen release from aqueous sodium formate. Three supports were used: carbon black (Vulcan XC-72), nitrogen-doped carbon (NC), and graphitic carbon nitride (g-C₃N₄). Nitrogen-doped carbon was obtained by heating a melamine-carbon black mixture at 700 °C under nitrogen, while g-C₃N₄ was synthesized by heating urea at 500 °C in air. Pd catalysts were produced by reducing H₂PdCl₄ with NaBH₄ in trisodium citrate solution at 25 °C, yielding a 1 wt% Pd loading. The product was filtered, washed, and dried at 85 °C for 24 h, and selected samples were calcined at 250 °C for 3 h in air. Structural and compositional analyses included inductively coupled plasma-optical emission spectrometry (PerkinElmer 7300 DV) to determine Pd content, X-ray diffraction (Rigaku SmartLab SE, Cu K α , 2 θ = 2–100°) to assess crystallinity, and nitrogen physisorption (Micromeritics ASAP 2020) using BET and BJH models to measure surface area and pore volume. Pd dispersion was quantified by CO chemisorption (Micromeritics ASAP 2020C, 30 °C, pre-reduced at 100 °C for 0.5 h), and nanoparticle morphology was examined by aberration-corrected scanning transmission electron microscopy (Thermo Fisher Themis Z, 300 kV). Catalytic performance was tested in a 50 mL batch reactor containing 250 mg of catalyst and 10 mL of 1 M sodium formate at 65 °C under N₂ with stirring at 500 rpm for 2 h, where gas evolution was monitored by pressure change and analyzed using a micro-gas chromatograph. In-situ X-ray absorption spectroscopy was performed at the Stanford Synchrotron Radiation Lightsource beamline 4-1 to monitor Pd oxidation states during reaction using Pd K-edge XANES and EXAFS scans (24126–25238 eV, 0.5 × 4 mm beam). Catalyst reuse tests were carried out by recovering the solid after reaction, washing with deionized water, drying at 80 °C, and re-calcining at 180 or 250 °C for 3 h when required. All synthesis, characterization, and catalytic experiments were conducted under controlled temperature and atmospheric conditions to ensure reproducibility.

Measurements Taken:

- Pd oxidation state and local atomic structure characterized by in-situ X-ray Absorption Spectroscopy (XAS, SSRL beamline 4-1) with Pd K-edge XANES and EXAFS scans (24126–25238 eV, beam size 0.5 × 4 mm) under reaction conditions.

- Palladium loading (wt%) measured using Inductively Coupled Plasma-Optical Emission Spectroscopy (ICP-OES, PerkinElmer 7300 DV) to quantify Pd content on carbon supports.

Outcome Prediction Question: Palladium nanoparticles supported on carbon materials were assessed as catalysts for hydrogen release from aqueous sodium formate. Three supports- carbon black (Vulcan XC-72), nitrogen-doped carbon (NC), and graphitic carbon nitride (g-C₃N₄)- were employed, with NC synthesized by heating a melamine-carbon black mixture at 700 °C under N₂ and g-C₃N₄ prepared by urea pyrolysis at 500 °C in air. Pd catalysts (1 wt%) were obtained by reducing H₂PdCl₄ with NaBH₄ in trisodium citrate at 25 °C, followed by drying and optional calcination at 250 °C. Structural and chemical characterization included ICP-OES for Pd content, XRD for crystallinity, N₂ physisorption for surface area, CO chemisorption for Pd dispersion, and STEM for nanoparticle morphology. Catalytic performance was evaluated in a batch reactor (65 °C, 1 M sodium formate) by monitoring gas evolution and composition via micro-GC. In-situ XANES/EXAFS at the Pd K-edge tracked oxidation-state changes during reaction, and reuse tests examined catalyst stability following washing and re-calcination. What will in-situ XANES analysis reveal about the role of palladium oxide (PdO) as an active catalyst for formate dehydrogenation?

Ground Truth Answer: In-situ XANES experiments unambiguously demonstrate that PdO is rapidly reduced to metallic Pd and then forms Pd hydride upon exposure to a formate solution, showing that PdO does not play a direct role in the mechanism of H₂ formation.

Background Knowledge:

- Palladium nanoparticles on carbon supports (Pd/C) are effective for catalyzing hydrogen release from aqueous formate solutions but typically suffer from a gradual decrease of activity.

- Nitrogen doping of carbon supports is observed to enhance the rates of hydrogen release from aqueous formate solutions

Rubrics: The response must state that palladium oxide (PdO) does not play a direct role as the active catalyst in the mechanism of H₂ formation.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Chemistry: Multiple-Choice Question

Paper Title: Lab-Scale Thermal Decomposition of Hydrogen Peroxide as Green Propellant over Low-Cost Catalysts Based on Copper Deposited on Different Supports

Link to The Paper: <https://www.mdpi.com/2226-4310/12/5/440>

Experimental Setup: Researchers investigate the thermal degradation of the H_2O_2 green monopropellant. Three distinct catalysts—copper supported on γ -alumina, graphite, and MNC clay—were used. Conversely, a LABSYS evo-gasorption apparatus (Category: DTA/T-G/DSC, Model: Setaram Instrumentation) was used to perform differential thermal analysis–thermogravimetry (DTA–TG) measurements in order to investigate the thermal breakdown of H_2O_2 at constant atmospheric pressure ($p = 1 \text{ atm}$). A syringe was used to inject a 30% (w/w) H_2O_2 microdroplet into the metallic sample cell. It was investigated how the three different catalysts affected the H_2O_2 thermogram. A microdroplet of liquid H_2O_2 was combined with a modest amount (a few micrograms) of powdered catalyst in the aluminum sample cell for each thermal study. Before each run, the following experimental conditions were maintained: (i) Carrier gas: argon, with a flow rate of $50 \text{ mL}\cdot\text{min}^{-1}$; (ii) Heating rate: $10 \text{ }^\circ\text{C}\cdot\text{min}^{-1}$, from room temperature up to $250 \text{ }^\circ\text{C}$; (iii) The H_2O_2 droplet was added directly to the catalyst particles already placed in the aluminum cell. After sealing the apparatus, a stabilization period of approximately 2 min was allowed for the system (carrier gas and sample) to equilibrate. The thermal run was then initiated to record the DTA–TG thermograms. Experiments were run at two constant temperatures: $0 \text{ }^\circ\text{C}$ and $36 \text{ }^\circ\text{C}$

Measurements Taken:

- Differential pressure (ΔP , in kPa) vs time (minutes) was recorded.
- ΔP for each catalyst (Cu/ γ -alumina, Cu/graphite, Cu/clay) compared to the uncatalyzed control.
- ΔP at $0 \text{ }^\circ\text{C}$ and $36 \text{ }^\circ\text{C}$ to assess temperature effects on decomposition rate.

Outcome Prediction Question: Which of the following statements best describes the observed catalytic activity (as measured by differential pressure, ΔP , vs time) for the decomposition of 30 % H_2O_2 over the three copper-supported catalysts (Cu/ γ -alumina, Cu/graphite, Cu/clay) compared to the uncatalyzed decomposition, at $36 \text{ }^\circ\text{C}$ and $0 \text{ }^\circ\text{C}$?

- A. At both temperatures all three catalysts produce rates almost identical to each other; the rates follow a similar trend, with 0°C just being slower than $36 \text{ }^\circ\text{C}$, each gives a large increase over the uncatalyzed reaction at both temperatures.
- B. At 0°C all three catalysts give a similar rate, none of them is clearly faster than another, but at $36 \text{ }^\circ\text{C}$ Cu/ γ -alumina gives the highest rate (largest ΔP increase), followed by Cu/graphite, then Cu/clay, each significantly faster than uncatalyzed at both temperatures.
- C. At $0 \text{ }^\circ\text{C}$, Cu/clay a rate that is slower than the uncatalyzed reaction at the beginning, then becomes faster than the uncatalyzed reaction, while Cu/graphite, and Cu/ γ -alumina have a similar rate and are higher than the uncatalyzed reaction. At $36 \text{ }^\circ\text{C}$ all three are faster than uncatalyzed reaction, Cu/ γ -alumina is the fastest, closely followed by Cu/graphite, then Cu/clay.
- D. At $0 \text{ }^\circ\text{C}$ all three catalysts begin slightly faster than the uncatalyzed reaction then all three become much faster, the variability being larger than the difference between the catalysts. At $36 \text{ }^\circ\text{C}$ the reaction with all three catalysts is much faster than the uncatalyzed reaction, with Cu/ γ -alumina being much faster than Cu/graphite, then Cu/clay lags because the copper particles came off the support particles.

Ground Truth Answer: C

Background Knowledge:

- As the world increasingly focuses on sustainable and environmentally friendly solutions, there is a growing interest in exploring greener alternative propellants that offer comparable performance while mitigating the drawbacks associated with hydrazine and its derivatives.
- The thermal decomposition of hydrogen peroxide (H_2O_2) as a promising green propellant was performed over free-noble metallic-based catalysts deposited on abundant supports.
- Green monopropellants have the potential for long-term cost savings due to reduced safety measures, disposal costs, and regulatory compliance requirements associated with hazardous materials such as hydrazine.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

Chemistry: Numerical Value Question

Paper Title: Time-resolved photo-electrochemical measurements to study band bending of BiVO₄ photoanodes

Link to The Paper: <https://chemrxiv.org/engage/chemrxiv/article-details/68b1a2e2728bf9025e19a17e?>

Experimental Setup: Thin-film BiVO₄ photoanodes were investigated in a three-electrode photo-electrochemical RRDE cell under chopped AM 1.5G illumination. Light switch-ON/OFF transients were recorded over 0–2.5 V vs RHE, and the disk photocurrent during switch-ON was fit with exponentials to isolate the fast space-charge reorganization time constant (τ_{fast}) (along with slower components).

Measurements Taken:

- Disk photocurrent transients at light switch-ON/OFF (current vs time) across 0–2.5 V vs RHE.
- Exponential fits of transients to extract characteristic time constants (including τ_{fast}) in seconds; report the average τ_{fast} (switch-ON) over the potential window.
- Steady-state J–E curves under illumination.
- RRDE ring current (Pt ring) vs time/potential for O₂ detection/validation.
- Assignment of τ_{fast} to space-charge reorganization based on transient behavior.

Outcome Prediction Question: Thin-film BiVO₄ photoanodes were tested in a three-electrode photo-electrochemical RRDE cell under chopped AM 1.5G illumination. During light “switch-ON” steps over 0–2.5 V vs RHE, the disk photocurrent transients were fit with exponentials to isolate the fast space-charge reorganization process (τ_{fast}). At these conditions, what is the average value of τ_{fast} in seconds (s) for the switch-ON process?

Ground Truth Answer: 0.0022±0.002 s.

Background Knowledge:

- BiVO₄ is a semiconductor photoanode used for oxygen evolution under illumination; its behavior is probed in a three-electrode photoelectrochemical cell.
- Band bending at the semiconductor/electrolyte interface creates a space-charge region that governs carrier separation and the early transient response.
- Time-resolved photoelectrochemistry with chopped AM 1.5G illumination measures photocurrent transients at light on/off to extract characteristic time constants.
- A rotating ring–disk electrode (RRDE) uses a Pt ring to detect dissolved O₂ produced at the disk, distinguishing disk photocurrent from ring current.
- The flat-band potential is the potential where band bending vanishes and is estimated from cyclic-voltammetry features; potentials are reported vs RHE.
- Exponential fitting of transients yields τ_{fast} and slower components that reflect interfacial charge reorganization and reaction kinetics.

1566 C.2 EXAMPLE HUMAN RESPONSES
15671568 Physics: Numerical Value Question
15691570 **Paper Title:** Recent Highlights from the STAR Experiment1571 **Link to The Paper:** <https://arxiv.org/abs/2508.08444>
15721573 **Experimental Setup:** Researchers investigated the Beam Energy Scan-II (BES-II) program
1574 at the STAR experiment, which was used to measure net-proton cumulant ratios in
1575 Gold-on-Gold (Au+Au) collisions at various center-of-mass energies (from 7.7 to 27
1576 GeV) in the Fixed-Target mode. BES-II employed a new centrality definition, RefMult3X,
1577 corresponding to pseudorapidity acceptances fulfilling $|\eta| < 1.6$. The Time-Projection
1578 Chamber (TPC) for low transverse momentum ($0.4 < p_T < 0.8$ GeV/c) and the Time-Of-Flight
1579 (TOF) detector for greater transverse momentum ($0.8 < p_T < 2.0$ GeV/c) were used to identify
1580 protons and anti-protons. Only particles falling within the speed window of $|y| < 0.5$ were
1581 included in the analysis. The most central collisions (0-5% centrality class) were the focus
1582 of the measurements, which were methodically adjusted for experimental variables such
1583 detector efficiency, event pile-up, and centrality bin width.1582 **Measurements Taken:**

- 1583 - Net-proton cumulants (C1, C2, C3, C4) as a function of collision centrality and collision
-
- 1584 energy.
-
- 1585 - The relative dynamical correlation of transverse momentum as a function of collision energy.

1586 **Outcome Prediction Question:** In the STAR experiment's Beam Energy Scan-II (BES-II),
1587 what was the measured value of the net-proton cumulant ratio C4/C2 at the collision energy
1588 of 19.6 GeV for the 0-5% centrality class?1589 **Ground Truth Answer:** [0.25-0.40]1590 Note: The range is informed graphically in Figure 3. The range was estimated by the pixel
1591 coordinates of the error bars and axis ticks.1592 **Background Knowledge:**

- 1593 - The upgrades done to STAR for BES-II enabled a new centrality definition, RefMult3X,
-
- 1594 which achieves better centrality resolution due to larger multiplicity within the acceptance.
-
- 1595 - Experimentally measured proton multiplicity distributions are described by the central
-
- 1596 moments, which depend on the cumulants. In particular, the second cumulant C2 is the
-
- 1597 variance
- σ^2
- , and the ratio between the fourth and second cumulant, C4/C2, is
- $\kappa\sigma^2$
- , where
- κ
- is
-
- 1598 the kurtosis.
-
- 1599 - When there are no intrinsic correlations among the measured particles, all ratios of the
-
- 1600 cumulants are unity, so Poisson statistics is a trivial baseline for experimentally measured
-
- 1601 cumulant ratios.
-
- 1602
-
- 1603
-
- 1604
-
- 1605
-
- 1606
-
- 1607
-
- 1608
-
- 1609
-
- 1610
-
- 1611
-
- 1612
-
- 1613
-
- 1614
-
- 1615
-
- 1616
-
- 1617
-
- 1618
-
- 1619

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Human Responses

Answer (NBK): The measured value of the net-proton cumulant ratio C_4/C_2 at the collision energy of 19.6 GeV for the 0-5% centrality class is 0.4.

Reasoning (NBK): Theoretically, the ratio of the fourth- to the second-order net-proton cumulant (C_4/C_2) is often called the moment product $\kappa\sigma^2$. Here, κ is the kurtosis and σ^2 is the variance. The theoretical Poisson baseline for net-proton cumulant ratios is unity or 1. Hence, the measured value must be ≤ 1 .

Additionally, various previous experiments support the fact that the C_4/C_2 value is close to unity for all collision energies for the smallest rapidity acceptance, and for higher collision energies. For example, according to Adam et al. (2021), during the BES-I experiment of the STAR detector at RHIC, the mean C_4/C_2 ratio in the 0-5% bin is ~ 0.4 . This result is also supported by Bleicher et. al. (1999) during the Ultra-Relativistic Quantum Molecular Dynamics (UrQMD) experiment.

Confidence (NBK): Somewhat confident in your answer

Difficulty (NBK): Easy to answer

Answer (BK): The measured value of the net-proton cumulant ratio C_4/C_2 at the collision energy of 19.6 GeV for the 0-5% centrality class is 0.4.

Reasoning (BK): Theoretically, the ratio of the fourth- to the second-order net-proton cumulant (C_4/C_2) is often called the moment product $\kappa\sigma^2$. Here, κ is the kurtosis and σ^2 is the variance. The theoretical Poisson baseline for net-proton cumulant ratios is unity or 1. Hence, the measured value must be ≤ 1 .

Additionally, various previous experiments support the fact that the C_4/C_2 value is close to unity for all collision energies for the smallest rapidity acceptance, and for higher collision energies. For example, according to Adam et al. (2021), during the BES-I experiment of the STAR detector at RHIC, the mean C_4/C_2 ratio in the 0-5% bin is ~ 0.4 [Figure 8]. This result is also supported by Bleicher et. al. (1999) during the Ultra-Relativistic Quantum Molecular Dynamics (UrQMD) experiment [Figures 6 and 30].

Confidence (BK): Somewhat confident in your answer

Difficulty (BK): Easy to answer

Feasibility: Very feasible to answer without running the experiment

Feasibility Reasoning: Theoretically, the ratio of the fourth- to the second-order net-proton cumulant (C_4/C_2) is often called the moment product $\kappa\sigma^2$. Here, κ is the kurtosis and σ^2 is the variance. The theoretical Poisson baseline for net-proton cumulant ratios is unity or 1. Hence, the measured value must be ≤ 1 , which can be directly concluded from the known theory on this topic.

Additionally, various previous experiments support the fact that the C_4/C_2 value is close to unity for all collision energies for the smallest rapidity acceptance, and for higher collision energies. For example, according to Adam et al. (2021), during the BES-I experiment of the STAR detector at RHIC, the mean C_4/C_2 ratio in the 0-5% bin is ~ 0.4 [Figure 8]. This result is also supported by Bleicher et. al. (1999) during the Ultra-Relativistic Quantum Molecular Dynamics (UrQMD) experiment [Figures 6 and 30].

Hence, using the existing literature on previously performed experiments, the measured value of C_4/C_2 can be logically estimated for the BES-II experiment of the STAR detector at RHIC.

D ADDITIONAL RESULTS

Table 3: **Domain-wise performance across LLM families.** Different versions of Gemini, OpenAI, Claude, Llama, Qwen, and DeepSeek evaluated on *Chemistry, Biology, Physics, and All Domains*. **Conf.** := Confidence Score; **Diff.** := Difficulty Level; **Feas.** := Feasibility Score.

Model	Experimental Setup	Chemistry				Biology				Physics				All Domains				
		Accuracy (%)		Calibration (1-5)		Accuracy (%)		Calibration (1-5)		Accuracy (%)		Calibration (1-5)		Accuracy (%)		Calibration (1-5)		
		Conf.	Diff.	Feas.	Feas.													
Gemini 3-pro	NBK	26.14±5.40	4.37±0.01	3.56±0.03	3.55±0.04	24.47±1.24	4.38±0.02	3.44±0.01	3.55±0.09	26.00±0.00	4.56±0.02	3.39±0.01	3.71±0.10	25.27±1.92	4.42±0.01	3.46±0.01	3.59±0.06	
	BK	28.43±0.98	4.39±0.03	3.52±0.02	3.50±0.07	27.29±0.75	4.39±0.02	3.40±0.07	3.59±0.03	26.33±2.08	4.56±0.01	3.25±0.06	3.89±0.06	27.33±0.79	4.43±0.01	3.40±0.04	3.64±0.02	
	SBK	28.43±0.00	4.43±0.00	3.45±0.00	3.71±0.00	26.11±0.00	4.45±0.00	3.39±0.00	3.89±0.00	21.00±0.00	4.63±0.00	3.24±0.00	4.00±0.00	25.43±0.00	4.49±0.00	3.27±0.00	3.87±0.00	
	SABK	30.39±0.00	4.48±0.00	3.42±0.00	3.79±0.00	25.12±0.00	4.44±0.00	3.20±0.00	3.78±0.00	27.00±0.00	4.68±0.00	3.14±0.00	4.01±0.00	26.91±0.00	4.51±0.00	3.24±0.00	3.84±0.00	
	FBK	25.99±0.00	4.35±0.00	3.53±0.00	3.47±0.00	24.63±0.00	4.32±0.00	3.49±0.00	3.35±0.00	23.00±0.00	4.63±0.00	3.32±0.00	3.75±0.00	24.44±0.00	4.40±0.00	3.46±0.00	3.48±0.00	
Claude Opus 4.5	NBK	15.69±0.98	3.19±0.08	4.04±0.03	2.78±0.08	25.12±0.40	4.30±0.04	3.98±0.02	2.92±0.04	26.33±2.08	3.48±0.03	4.01±0.02	3.17±0.06	23.05±0.51	3.23±0.04	4.00±0.02	2.95±0.01	
	BK	22.86±1.50	3.20±0.01	4.03±0.01	2.85±0.06	27.09±0.85	3.38±0.02	3.92±0.05	3.00±0.01	32.33±2.08	3.51±0.04	3.95±0.04	3.14±0.03	27.33±0.75	3.37±0.02	3.95±0.03	3.00±0.01	
	SBK	17.65±0.00	3.27±0.00	4.01±0.00	2.86±0.00	27.00±0.00	3.46±0.00	3.81±0.00	3.12±0.00	29.00±0.00	3.51±0.00	3.90±0.00	3.23±0.00	25.14±0.00	3.43±0.00	3.88±0.00	3.08±0.00	
	SABK	18.63±0.00	3.40±0.00	3.91±0.00	3.01±0.00	26.60±0.00	3.48±0.00	3.83±0.00	3.17±0.00	32.00±0.00	3.64±0.00	3.50±0.00	3.41±0.00	25.93±0.00	3.50±0.00	3.83±0.00	3.19±0.00	
	FBK	18.63±0.00	3.19±0.00	4.02±0.00	2.80±0.00	26.60±0.00	3.28±0.00	3.97±0.00	2.95±0.00	32.00±0.00	3.39±0.00	3.97±0.00	3.12±0.00	25.93±0.00	3.28±0.00	3.98±0.00	2.95±0.00	
Claude Sonnet 4.5	NBK	23.86±1.50	3.98±0.02	3.77±0.01	3.24±0.05	22.66±0.85	4.04±0.02	3.63±0.02	3.47±0.03	21.00±2.00	4.14±0.04	3.74±0.02	3.39±0.03	22.55±0.75	4.05±0.01	3.69±0.01	3.40±0.01	
	BK	26.80±1.50	4.05±0.03	3.75±0.03	3.31±0.02	28.08±2.15	4.06±0.02	3.57±0.03	3.54±0.03	24.67±1.53	4.10±0.03	3.67±0.02	3.47±0.06	26.91±1.23	4.07±0.02	3.64±0.02	3.47±0.02	
	SBK	18.10±0.00	4.10±0.00	3.73±0.00	3.53±0.00	20.00±0.00	4.11±0.00	3.57±0.00	3.64±0.00	20.00±0.00	4.11±0.00	3.79±0.01	3.52±0.00	19.16±0.00	4.11±0.00	3.64±0.00	3.58±0.00	
	SABK	19.61±0.00	4.11±0.00	3.75±0.00	3.42±0.00	25.12±0.00	4.11±0.00	3.52±0.00	3.66±0.00	29.00±0.00	4.06±0.00	3.82±0.00	3.47±0.00	24.99±0.00	4.10±0.00	3.65±0.00	3.55±0.00	
	FBK	23.53±0.00	3.99±0.00	3.80±0.00	3.16±0.00	23.65±0.00	3.94±0.00	3.67±0.00	3.40±0.00	22.00±0.00	3.95±0.00	3.85±0.00	3.33±0.00	23.21±0.00	3.96±0.00	3.75±0.00	3.32±0.00	
Claude Opus 4.1	NBK	22.55±2.59	4.00±0.06	3.76±0.03	2.99±0.06	20.36±0.28	4.01±0.01	3.61±0.01	3.20±0.00	25.67±4.16	4.09±0.03	3.77±0.01	3.29±0.02	22.27±1.48	4.03±0.01	3.69±0.01	3.17±0.01	
	BK	22.55±0.00	4.02±0.03	3.74±0.02	3.10±0.07	26.11±1.71	4.05±0.01	3.54±0.02	3.28±0.04	27.00±1.00	4.11±0.01	3.66±0.03	3.37±0.02	24.43±0.65	4.05±0.01	3.62±0.02	3.26±0.02	
	SBK	24.51±0.00	4.10±0.00	3.69±0.00	3.29±0.00	20.81±0.00	4.15±0.00	3.48±0.00	3.50±0.00	24.00±0.00	4.16±0.00	3.67±0.00	3.57±0.00	22.53±0.00	4.14±0.00	3.58±0.00	3.46±0.00	
	SABK	28.43±0.00	4.09±0.00	3.74±0.00	3.27±0.00	22.17±0.00	4.13±0.00	3.47±0.00	3.55±0.00	25.00±0.00	4.13±0.00	3.65±0.00	3.64±0.00	24.44±0.00	4.12±0.00	3.58±0.00	3.50±0.00	
	FBK	20.59±0.00	3.97±0.00	3.76±0.00	2.96±0.00	22.17±0.00	3.97±0.00	3.67±0.00	3.15±0.00	23.00±0.00	3.98±0.00	3.78±0.00	3.39±0.00	21.98±0.00	3.97±0.00	3.72±0.00	3.16±0.00	
Gemini 3-Flash	NBK	22.88±1.50	4.43±0.03	3.58±0.01	4.23±0.01	22.33±2.48	4.37±0.01	3.32±0.04	4.32±0.02	21.33±2.31	4.47±0.03	3.45±0.03	4.34±0.01	22.22±1.08	4.41±0.02	3.42±0.02	4.30±0.01	
	BK	24.84±4.08	4.42±0.03	3.56±0.04	4.24±0.01	23.97±1.24	4.41±0.01	3.25±0.02	4.35±0.02	23.00±1.00	4.52±0.04	3.39±0.08	4.37±0.05	23.95±1.62	4.44±0.01	3.36±0.02	4.33±0.01	
	SBK	23.53±0.00	4.50±0.00	3.48±0.00	4.26±0.00	21.78±0.00	4.47±0.00	3.19±0.00	4.41±0.00	20.83±0.00	4.51±0.00	3.38±0.00	4.36±0.00	21.99±0.00	4.49±0.00	3.31±0.00	4.36±0.00	
	SABK	24.48±0.00	4.48±0.00	3.45±0.00	4.26±0.00	20.00±0.00	4.45±0.00	3.50±0.00	4.30±0.00	26.33±0.00	4.51±0.00	3.49±0.00	4.39±0.00	26.37±0.00	4.51±0.00	3.39±0.00	4.39±0.00	
	FBK	29.41±0.00	4.40±0.00	3.58±0.00	4.21±0.00	23.65±0.00	4.34±0.00	3.37±0.00	4.28±0.00	21.00±0.00	4.44±0.00	3.51±0.00	4.29±0.00	24.44±0.00	4.38±0.00	3.46±0.00	4.26±0.00	
OpenAI GPT-5	NBK	18.95±2.04	3.53±0.04	3.49±0.01	3.61±0.05	20.69±1.48	3.59±0.02	3.37±0.03	3.60±0.02	22.00±1.73	3.61±0.02	3.37±0.03	3.67±0.05	20.58±1.03	3.58±0.02	3.40±0.03	3.62±0.03	
	BK	19.93±0.57	3.59±0.06	3.43±0.01	3.69±0.01	18.72±0.00	3.70±0.00	3.20±0.00	3.79±0.00	21.00±0.00	3.60±0.00	3.32±0.00	3.71±0.00	19.01±0.00	3.65±0.00	3.27±0.00	3.76±0.00	
	SBK	17.65±0.00	3.60±0.00	3.36±0.00	3.76±0.00	18.72±0.00	3.70±0.00	3.20±0.00	3.79±0.00	21.00±0.00	3.60±0.00	3.32±0.00	3.71±0.00	19.01±0.00	3.65±0.00	3.27±0.00	3.76±0.00	
	SABK	18.63±0.00	3.62±0.00	3.34±0.00	3.76±0.00	24.63±0.00	3.79±0.00	3.17±0.00	3.86±0.00	24.00±0.00	3.55±0.00	3.28±0.00	3.72±0.00	22.96±0.00	3.69±0.00	3.24±0.00	3.80±0.00	
	FBK	20.59±0.00	3.49±0.00	3.50±0.00	3.58±0.00	19.70±0.00	3.60±0.00	3.37±0.00	3.56±0.00	18.00±0.00	3.64±0.00	3.42±0.00	3.68±0.00	19.51±0.00	3.58±0.00	3.41±0.00	3.60±0.00	
Human Baseline	NBK	8.82	2.59	3.78	2.52	23.15	3.13	3.39	2.92	26.00	3.07	3.26	3.05	20.25	2.98	3.46	2.85	
	BK	2.65	3.78	2.52	23.65	2.92	3.32	2.92	27.00	2.92	3.05	3.10	3.42	2.85				
	NBK	19.28±5.91	4.34±0.02	3.25±0.06	4.41±0.01	21.02±1.50	4.38±0.02	3.08±0.04	4.46±0.03	18.00±2.00	4.42±0.05	3.16±0.03	4.46±0.01	19.84±1.49	4.38±0.01	3.14±0.03	4.44±0.02	
	BK	21.24±2.26	4.42±0.03	3.16±0.05	4.48±0.05	22.00±1.03	4.48±0.01	2.94±0.02	4.56±0.00	20.33±0.58	4.51±0.03	3.07±0.04	4.53±0.03	21.40±0.87	4.47±0.01	3.02±0.03	4.53±0.02	
	SBK	11.79±0.00	4.45±0.00	3.11±0.00	4.35±0.00	16.08±0.00	4.46±0.00	2.87±0.00	4.57±0.00	16.00±0.00	4.49±0.00	3.16±0.00	4.51±0.00	14.96±0.00	4.45±0.00	3.00±0.00	4.52±0.00	
OpenAI O3-mini	NBK	19.61±0.00	4.33±0.00	3.12±0.00	4.43±0.00	17.24±0.00	4.48±0.00	2.85±0.00	4.56±0.00	26.00±0.00	4.52±0.00	3.00±0.00	4.58±0.00	20.00±0.00	4.45±0.00	2.95±0.00	4.53±0.00	
	BK	17.65±0.00	4.29±0.00	3.25±0.00	4.39±0.00	19.70±0.00	4.36±0.00	3.11±0.00	4.45±0.00	19.00±0.00	4.42±0.00	3.20±0.00	4.44±0.00	19.01±0.00	4.46±0.00	3.17±0.00	4.43±0.00	
	BK	16.99±2.47	4.35±0.07	3.64±0.17	4.34±0.09	19.21±1.71	4.47±0.02	3.52±0.01	4.37±0.06	21.33±2.52	4.63±0.03	3.74±0.10	4.33±0.04	19.18±0.79	4.48±0.03	3.60±0.06	4.35±0.04	
	SBK	18.63±0.00	4.39±0.00	3.43±0.00	4.30±0.00	19.21±0.00	4.45±0.00	3.36±0.00	4.44±0.00	23.67±1.53	4.64±0.05	3.54±0.06	4.35±0.08	21.98±2.36	4.58±0.00	3.45±0.01	4.37±0.06	
	SBK	13.73±0.00	4.51±0.00	3.71±0.00	4.28±0.00	18.32±0.00	4.40±0.00	3.51±0.00	4.39±0.00	22.22±0.00	4.59±0.00	3.76±0.00	4.44±0.00	18.12±0.00	4.58±0.00	3.62±0.00	4.37±0.00	
DeepSeek v3	NBK	16.67±0.00	4.31±0.00	3.79±0.00	4.34±0.00	21.67±0.00	4.49±0.00	3.39±0.00	4.48±0.00	23.00±0.00	4.47±0.00	3.79±0.00	4.40±0.00	20.74±0.00	4.44±0.00	3.59±0.00	4.40±0.00	
	BK	16.99±2.26	3.53±0.03	3.65±0.03	3.44±0.07	19.54±0.28	3.51±0.05	3.68±0.02	3.38±0.00	16.67±0.58	3.47±0.02	3.72±0.04	3.31±0.03	18.19±0.71	3.50±0.03	3.68±0.01	3.38±0.01	
	BK	21.57±2.59	3.54±0.03	3.61±0.04	3.53±0.05	19.87±1.03	3.65±0.03	3.62±0.01	3.55±0.03	18.00±1.00	3.53±0.09	3.60±0.01	3.42±0.08	19.84±0.29	3.59±0.02	3.61±0.02	3.51±0.01	
	SBK	14.14±0.00	3.64±0.00	3.57±0.00	3.54±0.00	18.59±0.00	3.63±0.00	3.62±0.00	3.43±0.00	14.43±0.00	3.47±0.00	3.73±0.00	3.37±0.00	16.44±0.00	3.59±0.00	3.63±0.00	3.44±0.00	
	SABK	18.63±0.00	3.59±0.00	3.45±0.00	3.60±0.00	19.21±0.00	3.71±0.00	3.53±0.00	3.60±0.00	19.00±0.00	3.62±0.00	3.60±0.00	3.56±0.00	26.27±0.00	3.66±0.00	3.39±0.00	3.59±0.00	

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

Table 4: Question-format performance across LLM families. Different versions of Gemini, OpenAI, Claude (Opus/Sonnet), Llama, Qwen, and DeepSeek evaluated across question formats. **Conf.** := Confidence Score; **Diff.** := Difficulty Level; **Feas.** := Feasibility Score.

Model	Experimental Setup	MCQ						Numerical						Free form					
		Accuracy (%)		Calibration (1-5)			Accuracy (%)		Calibration (1-5)			Accuracy (%)		Calibration (1-5)					
		Conf.	Diff.	Feas.	Conf.	Diff.	Feas.	Conf.	Diff.	Feas.	Partial	Full	Conf.	Diff.	Feas.				
Gemini 3-pro	NBK	36.21 ± 1.78	4.42 ± 0.03	3.45 ± 0.01	3.89 ± 0.06	12.80 ± 2.58	4.19 ± 0.05	4.01 ± 0.03	2.46 ± 0.11	37.70 ± 1.65	22.39 ± 2.45	4.62 ± 0.01	3.01 ± 0.03	4.19 ± 0.09					
	BK	42.39 ± 1.55	4.46 ± 0.01	3.36 ± 0.05	3.96 ± 0.02	12.80 ± 1.36	4.16 ± 0.03	4.03 ± 0.02	2.47 ± 0.06	36.04 ± 0.78	21.12 ± 2.33	4.62 ± 0.03	2.90 ± 0.05	4.25 ± 0.03					
	SBK	37.04 ± 0.00	4.52 ± 0.00	3.17 ± 0.00	4.15 ± 0.00	10.71 ± 0.00	4.19 ± 0.00	3.98 ± 0.00	2.77 ± 0.00	38.90 ± 0.00	23.66 ± 0.00	4.71 ± 0.00	2.78 ± 0.00	4.47 ± 0.00					
	FBK	38.89 ± 0.00	4.41 ± 0.00	3.44 ± 0.00	3.81 ± 0.00	8.93 ± 0.00	4.14 ± 0.00	4.04 ± 0.00	2.28 ± 0.00	36.41 ± 0.00	19.85 ± 0.00	4.62 ± 0.00	2.98 ± 0.00	4.10 ± 0.00					
Claude Opus 4.5	NBK	33.54 ± 0.71	3.78 ± 0.05	3.88 ± 0.02	3.16 ± 0.01	13.99 ± 0.52	2.18 ± 0.03	4.59 ± 0.04	2.00 ± 0.02	34.69 ± 0.43	17.81 ± 1.17	3.74 ± 0.06	3.64 ± 0.01	3.50 ± 0.03					
	BK	39.09 ± 0.94	3.90 ± 0.03	3.80 ± 0.05	3.27 ± 0.02	15.77 ± 2.73	2.11 ± 0.02	4.58 ± 0.01	1.98 ± 0.02	38.37 ± 0.62	22.65 ± 1.17	3.79 ± 0.01	3.60 ± 0.03	3.54 ± 0.02					
	SBK	36.88 ± 0.00	3.95 ± 0.00	3.77 ± 0.00	3.34 ± 0.00	15.18 ± 0.00	2.23 ± 0.00	4.52 ± 0.00	2.08 ± 0.00	35.13 ± 0.00	19.23 ± 0.00	3.81 ± 0.00	3.47 ± 0.00	3.61 ± 0.00					
	FBK	35.80 ± 0.00	3.98 ± 0.00	3.69 ± 0.00	3.54 ± 0.00	17.86 ± 0.00	2.33 ± 0.00	4.53 ± 0.00	2.15 ± 0.00	36.34 ± 0.00	20.61 ± 0.00	3.91 ± 0.00	3.42 ± 0.00	3.64 ± 0.00					
Claude Sonnet 4.5	NBK	29.01 ± 1.63	4.22 ± 0.03	3.60 ± 0.05	3.65 ± 0.01	16.07 ± 0.89	3.59 ± 0.07	4.15 ± 0.05	1.99 ± 0.01	35.75 ± 2.22	20.10 ± 1.17	4.23 ± 0.01	3.42 ± 0.01	3.94 ± 0.05					
	BK	25.61 ± 1.28	4.26 ± 0.02	3.55 ± 0.05	3.80 ± 0.04	15.77 ± 2.58	3.56 ± 0.04	4.12 ± 0.02	2.33 ± 0.05	40.83 ± 0.40	25.19 ± 3.05	4.26 ± 0.01	3.35 ± 0.00	4.01 ± 0.04					
	SBK	25.62 ± 0.00	4.29 ± 0.00	3.50 ± 0.00	3.95 ± 0.00	11.93 ± 0.00	3.70 ± 0.00	4.11 ± 0.00	2.61 ± 0.00	36.26 ± 0.00	17.19 ± 0.00	4.23 ± 0.00	3.44 ± 0.00	3.95 ± 0.00					
	FBK	33.33 ± 0.00	4.27 ± 0.00	3.28 ± 0.00	4.37 ± 0.00	9.32 ± 0.00	4.22 ± 0.00	3.89 ± 0.00	4.05 ± 0.00	40.93 ± 0.00	25.95 ± 0.00	4.77 ± 0.00	2.79 ± 0.00	4.69 ± 0.00					
Claude Opus 4.1	NBK	29.01 ± 2.83	4.15 ± 0.02	3.64 ± 0.02	3.38 ± 0.06	16.07 ± 1.79	3.75 ± 0.07	4.00 ± 0.03	2.28 ± 0.07	34.38 ± 0.37	19.08 ± 0.76	4.13 ± 0.01	3.47 ± 0.01	3.68 ± 0.01					
	BK	35.39 ± 0.36	4.18 ± 0.01	3.54 ± 0.00	3.55 ± 0.04	14.58 ± 2.25	3.79 ± 0.05	4.01 ± 0.03	2.25 ± 0.05	37.52 ± 1.14	22.39 ± 0.44	4.12 ± 0.03	3.38 ± 0.03	3.76 ± 0.01					
	SBK	28.75 ± 0.00	4.20 ± 0.00	3.52 ± 0.00	3.76 ± 0.00	17.27 ± 0.00	3.95 ± 0.00	3.96 ± 0.00	2.43 ± 0.00	34.63 ± 0.00	19.38 ± 0.00	4.22 ± 0.00	3.33 ± 0.00	3.98 ± 0.00					
	FBK	30.86 ± 0.00	4.22 ± 0.00	3.52 ± 0.00	3.79 ± 0.00	16.96 ± 0.00	3.89 ± 0.00	3.99 ± 0.00	2.45 ± 0.00	38.24 ± 0.00	22.90 ± 0.00	4.18 ± 0.00	3.31 ± 0.00	4.03 ± 0.00					
Gemini 3-Flash	NBK	29.22 ± 2.85	4.38 ± 0.02	3.38 ± 0.01	4.30 ± 0.02	11.90 ± 1.03	4.18 ± 0.01	3.92 ± 0.03	4.02 ± 0.02	35.71 ± 1.56	22.39 ± 1.17	4.63 ± 0.02	3.03 ± 0.06	4.55 ± 0.04					
	BK	35.39 ± 1.28	4.44 ± 0.01	3.35 ± 0.02	4.32 ± 0.01	9.52 ± 1.86	4.18 ± 0.02	3.89 ± 0.02	4.03 ± 0.01	37.11 ± 3.11	22.14 ± 2.02	4.68 ± 0.03	2.92 ± 0.05	4.60 ± 0.02					
	SBK	31.06 ± 0.00	4.50 ± 0.00	3.24 ± 0.00	4.35 ± 0.00	12.61 ± 0.00	4.20 ± 0.00	3.89 ± 0.00	4.04 ± 0.00	33.56 ± 0.00	18.75 ± 0.00	4.72 ± 0.00	2.90 ± 0.00	4.65 ± 0.00					
	FBK	33.33 ± 0.00	4.33 ± 0.00	3.43 ± 0.00	4.26 ± 0.00	14.29 ± 0.00	4.21 ± 0.00	3.89 ± 0.00	3.98 ± 0.00	36.75 ± 0.00	22.14 ± 0.00	4.59 ± 0.00	3.11 ± 0.00	4.50 ± 0.00					
OpenAI GPT-5.2	NBK	28.81 ± 3.40	3.95 ± 0.00	3.12 ± 0.03	3.95 ± 0.02	12.50 ± 2.36	2.77 ± 0.05	4.02 ± 0.01	2.84 ± 0.06	33.80 ± 1.00	17.30 ± 1.17	3.82 ± 0.02	3.21 ± 0.04	3.88 ± 0.02					
	BK	30.04 ± 2.49	3.98 ± 0.02	3.01 ± 0.01	4.01 ± 0.02	13.10 ± 3.72	2.94 ± 0.07	3.99 ± 0.02	3.04 ± 0.05	38.52 ± 1.46	22.14 ± 1.32	3.92 ± 0.03	3.11 ± 0.02	3.94 ± 0.03					
	SBK	26.54 ± 0.00	3.92 ± 0.00	2.96 ± 0.00	4.01 ± 0.00	13.39 ± 0.00	2.97 ± 0.00	3.99 ± 0.00	3.13 ± 0.00	29.33 ± 0.00	14.50 ± 0.00	3.90 ± 0.00	3.03 ± 0.00	4.00 ± 0.00					
	FBK	27.16 ± 0.00	3.97 ± 0.00	2.92 ± 0.00	4.03 ± 0.00	16.96 ± 0.00	3.00 ± 0.00	3.94 ± 0.00	3.24 ± 0.00	36.06 ± 0.00	22.90 ± 0.00	3.92 ± 0.00	3.04 ± 0.00	3.99 ± 0.00					
Human Baseline	NBK	26.54	3.33	3.22	3.01	8.93	2.29	3.99	2.39	36.09	22.14	3.13	3.31	3.05					
	BK	27.16	3.46	3.14	3.01	9.82	2.36	3.96	2.39	36.86	22.90	3.28	3.30	3.05					
	SBK	26.45 ± 1.23	4.56 ± 0.02	2.76 ± 0.06	4.69 ± 0.02	14.58 ± 1.86	3.95 ± 0.03	3.73 ± 0.02	4.01 ± 0.03	29.95 ± 1.12	16.03 ± 2.75	4.51 ± 0.04	3.10 ± 0.02	4.51 ± 0.05					
	FBK	27.16 ± 0.00	4.69 ± 0.00	2.61 ± 0.00	4.77 ± 0.00	11.32 ± 0.00	3.95 ± 0.00	3.64 ± 0.00	4.07 ± 0.00	28.35 ± 0.00	14.06 ± 0.00	4.61 ± 0.00	2.96 ± 0.00	4.60 ± 0.00					
OpenAI O3-mini	NBK	23.66 ± 1.43	4.70 ± 0.04	3.50 ± 0.01	4.37 ± 0.04	13.39 ± 0.00	3.85 ± 0.13	4.08 ± 0.10	4.17 ± 0.10	33.14 ± 0.81	18.58 ± 1.92	4.74 ± 0.02	3.34 ± 0.11	4.49 ± 0.03					
	BK	27.98 ± 4.20	4.72 ± 0.01	3.29 ± 0.07	4.45 ± 0.04	12.50 ± 2.36	3.98 ± 0.04	3.97 ± 0.17	4.18 ± 0.08	36.92 ± 1.16	22.65 ± 2.68	4.84 ± 0.02	3.20 ± 0.11	4.43 ± 0.08					
	SBK	25.00 ± 0.00	4.75 ± 0.00	3.41 ± 0.00	4.44 ± 0.00	8.04 ± 0.00	4.10 ± 0.00	4.14 ± 0.00	4.11 ± 0.00	34.21 ± 0.00	18.32 ± 0.00	4.76 ± 0.00	3.44 ± 0.00	4.51 ± 0.00					
	FBK	25.31 ± 0.00	4.73 ± 0.00	3.36 ± 0.00	4.46 ± 0.00	10.71 ± 0.00	3.96 ± 0.00	4.05 ± 0.01	4.13 ± 0.00	36.37 ± 0.00	21.37 ± 0.00	4.77 ± 0.00	3.32 ± 0.00	4.52 ± 0.00					
DeepSeek v3	NBK	26.75 ± 0.71	3.92 ± 0.03	3.42 ± 0.02	3.83 ± 0.05	12.50 ± 0.89	2.73 ± 0.05	3.99 ± 0.02	2.67 ± 0.08	25.61 ± 1.94	12.47 ± 1.17	3.63 ± 0.08	3.74 ± 0.03	3.41 ± 0.01					
	BK	28.81 ± 0.26	4.01 ± 0.04	3.33 ± 0.02	3.92 ± 0.10	14.29 ± 1.79	2.84 ± 0.04	3.98 ± 0.03	2.82 ± 0.07	26.58 ± 1.46	13.49 ± 1.17	3.72 ± 0.06	3.64 ± 0.03	3.60 ± 0.05					
	SBK	18.24 ± 0.00	4.69 ± 0.00	2.61 ± 0.00	4.77 ± 0.00	11.32 ± 0.00	3.95 ± 0.00	3.64 ± 0.00	4.07 ± 0.00	28.35 ± 0.00	14.06 ± 0.00	4.61 ± 0.00	2.96 ± 0.00	4.60 ± 0.00					
	FBK	27.78 ± 0.00	4.69 ± 0.00	2.53 ± 0.00	4.78 ± 0.00	11.61 ± 0.00	3.92 ± 0.00	3.57 ± 0.00	4.02 ± 0.00	33.18 ± 0.00	17.56 ± 0.00	4.62 ± 0.00	2.95 ± 0.00	4.67 ± 0.00					
Llama 3.3 70B	NBK	23.66 ± 1.43	4.70 ± 0.04	3.50 ± 0.01	4.37 ± 0.04	13.39 ± 0.00	3.85 ± 0.13	4.08 ± 0.10	4.17 ± 0.10	33.14 ± 0.81	18.58 ± 1.92	4.74 ± 0.02	3.34 ± 0.11	4.49 ± 0.03					
	BK	27.98 ± 4.20	4.72 ± 0.01	3.29 ± 0.07	4.45 ± 0.04	12.50 ± 2.36	3.98 ± 0.04	3.97 ± 0.17	4.18 ± 0.08	36.92 ± 1.16	22.65 ± 2.68	4.84 ± 0.02	3.20 ± 0.11	4.43 ± 0.08					
	SBK	25.00 ± 0.00	4.75 ± 0.00	3.41 ± 0.00	4.44 ± 0.00	8.04 ± 0.00	4.10 ± 0.00	4.14 ± 0.00	4.11 ± 0.00	34.21 ± 0.00	18.32 ± 0.00	4.76 ± 0.00	3.44 ± 0.00	4.51 ± 0.00					
	FBK	25.31 ± 0.00	4.73 ± 0.00	3.36 ± 0.00	4.46 ± 0.00	10.71 ± 0.00	3.96 ± 0.00	4.05 ± 0.01	4.13 ± 0.00	36.37 ± 0.00	21.37 ± 0.00	4.77 ± 0.00	3.32 ± 0.00	4.52 ± 0.00					
OpenAI O3	NBK	26.75 ± 0.71	3.92 ± 0.03	3.42 ± 0.02	3.83 ± 0.05	12.50 ± 0.89	2.73 ± 0.05	3.99 ± 0.02	2.67 ± 0.08	25.61 ± 1.94	12.47 ± 1.17	3.63 ± 0.08	3.74 ± 0.03	3.41 ± 0.01					
	BK	28.81 ± 0.26	4.01 ± 0.04	3.33 ± 0.02	3.92 ± 0.10	14.29 ± 1.79	2.84 ± 0.04	3.98 ± 0.03	2.82 ± 0.07	26.58 ± 1.46	13.49 ± 1.17	3.72 ± 0.06	3.64 ± 0.03	3.60 ± 0.05					
	SBK	18.24 ± 0.00	4.69 ± 0.00	2.61 ± 0.00	4.77 ± 0.00	11.32 ± 0.00	3.95 ± 0.00	3.64 ± 0.00	4.07 ± 0.00	28.35 ± 0.00	14.06 ± 0.00	4.61 ± 0.00	2.96 ± 0.00	4.60 ± 0.00					
	FBK	27.78 ± 0.00	4.69 ± 0.00	2.53 ± 0.00	4.78 ± 0.00	11.61 ± 0.00	3.92 ± 0.00	3.57 ± 0.00	4.02 ± 0.00	33.18 ± 0.00	17.56 ± 0.00	4.62 ± 0.00	2.95 ± 0.00	4.67 ± 0.00					
OpenAI O1	NBK	26.75 ± 0.71	3.92 ± 0.03	3.42 ± 0.02	3.83 ± 0.05	12.50 ± 0.89	2.73 ± 0.05	3.99 ± 0.02	2.67 ± 0.08	25.61 ± 1.94	12.47 ± 1.17	3.63 ± 0.08	3.74 ± 0.03	3.41 ± 0.01					
	BK	29.22 ± 2.57	4.00 ± 0.01	2.96 ± 0.02	4.06 ± 0.01	12.80 ± 2.87	3.74 ± 0.08	3.19 ± 0.01	3.91 ± 0.01	37.99 ± 2.50	21.37 ± 2.75	3.99 ± 0.01	3.00 ± 0.01	4.03 ± 0.01					
	SBK	25.31 ± 0.00	4.02 ± 0.00	2.88 ± 0.00	4.12 ± 0.00	14.29 ± 0.00	3.80 ± 0.00	3.19 ± 0.00	3.94 ± 0.00	33.38 ± 0.00	17.56 ± 0.00	4.00 ± 0.00	2.95 ± 0.00	4.08 ± 0.00					
	FBK	24.07 ± 0.00	4.01 ± 0.00	2.90 ± 0.00	4.09 ± 0.00	14.29 ± 0.00	3.83 ± 0.00	3.19 ± 0.00	3.96 ± 0.00	35.89 ± 0.00	20.61 ± 0.00	4.00 ± 0.00	2.95 ± 0.00	4.10 ± 0.00					
Qwen 3 32B	NBK	22.43 ± 1.98	4.04 ± 0.01	3.82 ± 0.03	3.88 ± 0.09	12.50 ± 4.46	3.48 ± 0.03	4.09 ± 0.04	3.24 ± 0.01	28.35 ± 0.43	14.25 ± 1.59	3.97 ± 0.03	3.78 ± 0.05	3.86 ± 0.04					
	BK	23.87 ± 2.34	4.10 ± 0.03	3.75 ± 0.02	3.99 ± 0.02	13.39 ± 3.09	3.64 ± 0.01	4.01 ± 0.03	3.45 ± 0.04	31.07 ± 2.49	18.07 ± 2.68	4.01 ± 0.01	3.75 ± 0.04	3.85 ± 0.01					
	SBK	19.71 ± 0.00	4.06 ± 0.00	3.76 ± 0.00	3.94 ± 0.00	10.11 ± 0.00	3.66 ± 0.00	4.08 ± 0.00	3.51 ± 0.00	32.38 ± 0.00	17.65 ± 0.00	4.00 ± 0.00							

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

Table 7: Different versions of Gemini, OpenAI, Claude Sonnet, Llama, Qwen, and Deepseek evaluated on their ability to answer questions based on required background knowledge needed to answer questions.

Model	# Corr.	# Ques.	Acc (%)
Gemini 3-pro	1268	1350	93.93
Claude Opus 4.5	1277	1344	95.01
Claude Sonnet 4.5	1232	1316	93.62
Claude Opus 4.1	1228	1327	92.54
Gemini 3-Flash	1279	1350	94.74
OpenAI GPT-5.2	1276	1350	94.52
OpenAI O3-mini	1250	1350	92.59
DeepSeek v3	1234	1353	91.20
Llama 3.3 70B	1132	1350	83.85
OpenAI O3	1261	1350	93.41
Qwen 3 32B	1149	1342	85.62
Gemini 2.5-pro	1246	1350	92.30
Qwen 3 235B	1222	1350	90.52
OpenAI O4-mini	1252	1350	92.74
Llama 3.1 8B	955	1329	71.86

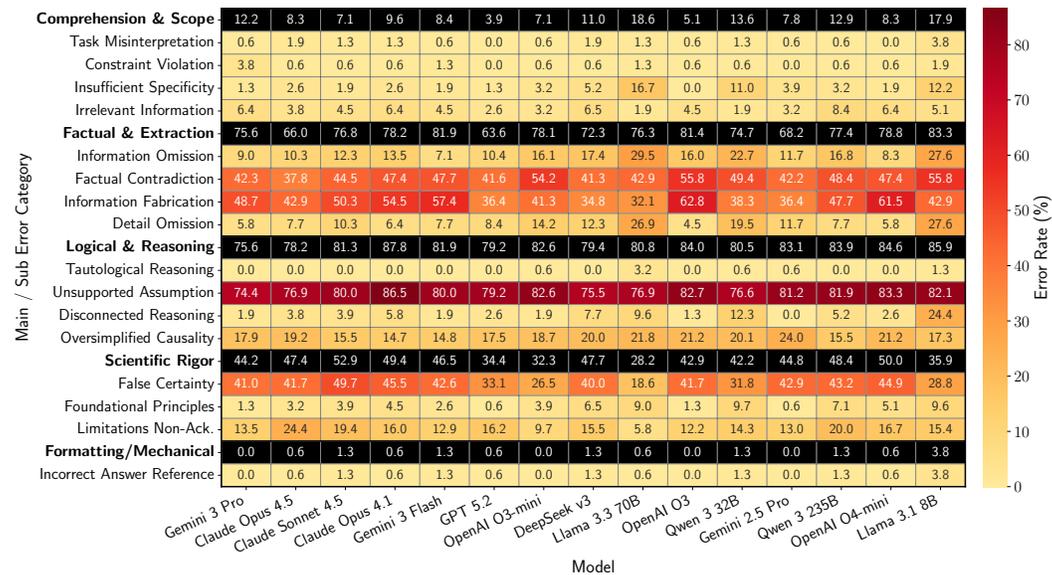


Figure 13: Analysis of model errors for high feasible questions. We employ an LLM judge to systematically classify errors in model predictions according to a hierarchical taxonomy spanning five top-level (in black background) categories and 16 specific error types. The heatmap shows the percentage of incorrect responses containing each error type for each evaluated model. Error categories progress from surface-level issues (Comprehension & Scope) to deeper reasoning failures (Logical & Reasoning Flaws) to fundamental scientific deficiencies (Deficiencies in Scientific Rigor). Models can exhibit multiple error types simultaneously, so accumulative percentage scores within top-level categories may exceed 100%. SciPredict tasks contribute to top-level category percentages if flagged with at least one underlying error type. Error analysis only considers the questions human experts marked as feasible to answer without running the practical experiment. Fig. 7 shows the same chart for all tasks. Tab. 8 provides comprehensive definitions for the error categories.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Table 8: Definitions of prediction error categories and flags.

Main Category	Specific Error	Description
Comprehension & Scope	Task Misinterpretation	The answer addresses a fundamentally different scientific question than the one that was asked.
	Constraint Violation	The answer ignores or violates a specific instruction or constraint mentioned in the question.
	Insufficient Specificity	The answer is overly generic, too high-level, or omits necessary details required to fully address the question.
	Irrelevant Information	The answer includes factually correct information that is not essential and does not help in answering the specific question.
Factual & Extraction	Information Omission	The answer fails to use a required piece of data that is explicitly present in the provided materials (e.g., experimental setup, measurements).
	Factual Contradiction	The answer misrepresents or directly contradicts facts, values, or relationships that are clearly stated in the provided materials.
	Information Fabrication	The answer invents data, formulas, or external “facts” that are not supported by the provided materials.
	Detail Omission	The answer’s own reasoning is incomplete because it omits a critical piece of evidence from the provided materials needed to support its conclusion.
Logical & Reasoning	Tautological Reasoning	The answer’s justification is circular, merely restating the conclusion in different words without providing independent evidence.
	Unsupported Assumption	The reasoning relies on a significant, unstated assumption that is not supported by the provided materials.
	Disconnected Reasoning	The answer lists correct facts but fails to logically connect them to form a coherent argument for the final conclusion.
	Oversimplified Causality	The reasoning focuses on a minor cause while ignoring a more critical or explicitly stated factor that impacts the conclusion.
Scientific Rigor	False Certainty	The answer presents a probabilistic, uncertain, or correlational outcome as a definitive fact, using absolute language where nuance is required.
	Foundational Principles	The reasoning violates a fundamental, universally accepted scientific principle to reach its conclusion.
	Limitations Non-Ack.	The answer presents a conclusion without acknowledging critical limitations, uncertainties, or sources of error evident from the experimental setup.
Formatting/Mechanical	Incorrect Answer Reference	Specific to multiple-choice questions, the reasoning correctly identifies one answer, but the final choice points to a different option letter.

E PROMPTS

Prompt used for errors analysis judge

[SYS]

Fields: domain, field

Instructions: You are acting as a judge evaluating a 'suggested_answer' to a scientific 'question' (of type 'question_type') which corresponds to the prediction of the outcome of a scientific experiment in {domain} and the field of {field}. Your goal is to identify the reason(s) why the provided answer is flawed or incorrect when compared to the 'ground_truth_answer' and the provided 'experimental_setup', 'measurements_taken', and 'background_knowledge'. Carefully review the provided materials and provide your judgment based on the rigorous definitions below. Your judgment should be based on a detailed analysis of the 'suggested_answer's reasoning and factual claims.

Evaluation Materials and Terminology:

- 'question': The scientific question posed to the responder for prediction of the experimental outcome.
- 'experimental_setup': Details of the experimental design, conditions, and procedures relevant to the 'question' provided to the responder for prediction of the experimental outcome.
- 'measurements_taken': Information about the measurements taken relevant to the 'question' provided to the responder for prediction of the experimental outcome.
- 'background_knowledge' (if any): Additional scientific context or principles relevant to the 'question' provided to the responder for prediction of the experimental outcome.
- 'suggested_answer': The responder's answer to the 'question', including any reasoning or justification provided.
- 'ground_truth_answer': The ground truth answer to the 'question', representing the correct prediction of the experimental outcome based on the provided materials.

Question Types:

- Multiple-Choice (MCQ): Includes a set of possible answers from which one (1) OR more (>1) must be selected.
- Free-Form: Requires a comprehensive but concise explanation of the expected experimental results.
- Numerical: Requires a specific numerical value prediction based on the provided data for the outcome of the experiment described in the question.

Error Analysis Categories:

1. Comprehension & Scope Errors: The answer fails because it fundamentally misunderstands the user's question or violates its core constraints. This is the primary error if the answer, regardless of its correctness, is for the wrong question.
2. Factual & Extraction Errors: The answer fails because it incorrectly handles explicit information from the provided 'experimental_setup', 'measurements_taken', or 'background_knowledge'. It omits, fabricates, or directly contradicts facts that are clearly stated.
3. Logical & Reasoning Flaws: The answer fails because the argument is logically unsound, even if the individual facts cited are correct. The connections between evidence and conclusion are invalid.
4. Deficiencies in Scientific Rigor: The answer fails because it lacks the necessary nuance and rigor expected in scientific communication. It may be factually correct but is presented with false certainty or violates a core scientific principle.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

5. Formatting & Mechanical Bug: The answer fails due to a non-substantive formatting error.

Detailed Analysis Flags:

First, choose a PRIMARY ERROR CATEGORY from the five main categories above that best explains WHY the 'suggested_answer' is flawed or incorrect. For this choice of the primary error category, provide a comprehensive justification (4-5 sentences) explaining your judgment.

Second, for EACH flag below (INCLUDING from ALL categories, NOT just the one you selected), choose YES, NO, or N/A based on the strict definitions provided:

1. Comprehension & Scope Errors

- 'flag_task_misinterpretation':

- Evidence Source: 'question', 'suggested_answer'.
- Definition: Whether the 'suggested_answer' addresses a fundamentally different question than the one posed.
- Prerequisite: None.
- 'YES': The answer's core purpose is different from the question's intent or it addresses a different scientific question than was asked.
- 'NO': The conditions for 'YES' are NOT satisfied.

- 'flag_constraint_violation':

- Evidence Source: 'question', 'suggested_answer'.
- Definition: Whether the 'suggested_answer' ignores a specific instruction or constraint mentioned in the 'question'.
- Prerequisite: The 'question' contains an explicit constraint.
- 'YES': The answer violates an explicit constraint in the question.
- 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
- 'N/A': The prerequisite is NOT met.

- 'flag_insufficient_specificity':

- Evidence Source: 'question', 'suggested_answer', 'ground_truth_answer'.
- Definition: Whether the 'suggested_answer' is overly generic or lacks the required detail.
- Prerequisite: None.
- 'YES': The answer is too high-level and omits details that are necessary to fully address the question, as evidenced by the 'ground_truth_answer'.
- 'NO': The conditions for 'YES' are NOT satisfied.

- 'flag_irrelevant_information':

- Evidence Source: 'question', 'suggested_answer', 'ground_truth_answer'.
- Definition: Whether the 'suggested_answer' includes factually correct but non-essential information.
- Prerequisite: None.
- 'YES': The answer contains information that does not help answer the specific 'question'.
- 'NO': The conditions for 'YES' are NOT satisfied.

2. Factual & Extraction Errors

- 'flag_information_omission':

- Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
- Definition: Whether the 'suggested_answer' fails to extract or reports as "missing" a REQUIRED piece of data explicitly present in the provided materials.
- Prerequisite: The information is explicitly stated in the 'experimental_setup', 'measurements_taken', or 'background_knowledge' AND the information is REQUIRED for answering the question.
- 'YES': A key fact, value, or condition from the provided materials is missing from, or was ignored in the 'suggested_answer'.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

- 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.
 - 'flag_factual_contradiction':
 - Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
 - Definition: Whether the 'suggested_answer' directly misrepresents or contradicts facts, values, or relationships stated in the provided materials.
 - Prerequisite: None.
 - 'YES': A statement in the 'suggested_answer' is verifiably FALSE when checked against the provided materials.
 - 'NO': The conditions for 'YES' are NOT satisfied.
 - 'flag_information_fabrication':
 - Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
 - Definition: Whether the 'suggested_answer' invents data, formulas, or external "facts" not supported by the provided materials.
 - Prerequisite: None.
 - 'YES': The answer includes specific information that cannot be found in or reasonably inferred from the provided materials.
 - 'NO': The conditions for 'YES' are NOT satisfied.
 - 'flag_detail_omission_in_reasoning':
 - Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
 - Definition: Whether the reasoning in the 'suggested_answer' omits a CRITICAL piece of evidence from the provided materials that is necessary to logically support its OWN conclusion.
 - Prerequisite: The 'suggested_answer' presents a logical argument or reasoning.
 - 'YES': The argument or reasoning provided for the answer is incomplete because a necessary premise from the provided materials is missing.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.
3. Logical & Reasoning Flaws
- 'flag_tautological_reasoning':
 - Evidence Source: 'suggested_answer'.
 - Definition: Whether the justification restates the conclusion without providing independent evidence.
 - Prerequisite: The 'suggested_answer' provides a justification or reasoning.
 - 'YES': The reasoning is circular, using the conclusion as its own evidence.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.
 - 'flag_unsupported_assumption':
 - Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
 - Definition: Whether the reasoning relies on a significant, unstated assumption that is NOT supported by the provided materials.
 - Prerequisite: The 'suggested_answer' presents a logical argument or reasoning.
 - 'YES': The logical leap from evidence to conclusion requires an assumption that is NOT provided or justified by the provided materials.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.
 - 'flag_disconnected_reasoning':
 - Evidence Source: 'suggested_answer'.
 - Definition: Whether the 'suggested_answer' lists correct facts but fails to logically connect them to the final conclusion.
 - Prerequisite: The 'suggested_answer' presents more than one (>1) piece of evidence in its reasoning.

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

- 'YES': NO logical connection is made between the evidence presented and the conclusion drawn.
- 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
- 'N/A': The prerequisite is NOT met.

- 'flag_oversimplified_causality':
 - Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
 - Definition: Whether the reasoning focuses on a minor cause while ignoring a more critical or explicitly stated factor impacting the conclusion to be made from the provided materials.
 - Prerequisite: The provided materials present multiple potential causal factors.
 - 'YES': The reasoning incorrectly prioritizes a secondary factor over the primary factor described in the provided materials.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.

- 4. Deficiencies in Scientific Rigor
 - 'flag_false_certainty':
 - Evidence Source: 'experimental_setup', 'measurements_taken', 'background_knowledge', 'suggested_answer'.
 - Definition: Whether the 'suggested_answer' presents a probabilistic, correlational, or uncertain outcome as a definitive fact.
 - Prerequisite: The outcome described in the provided materials or 'ground_truth_answer' is NON-deterministic.
 - 'YES': The answer uses absolute language where uncertainty or probability is warranted.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.

 - 'flag_violation_of_foundational_principles':
 - Evidence Source: 'suggested_answer'.
 - Definition: Whether the reasoning in the 'suggested_answer' is scientifically invalid because it violates a fundamental, universally accepted scientific principle.
 - Prerequisite: The 'suggested_answer' invokes reasoning related to a known scientific principle.
 - 'YES': The reasoning makes a statement that is verifiably FALSE according to a FOUNDATIONAL scientific principle.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.

 - 'flag_failure_to_acknowledge_limitations':
 - Evidence Source: 'experimental_setup', 'suggested_answer'.
 - Definition: Whether the 'suggested_answer' presents a conclusion without acknowledging critical limitations or uncertainties evident from the 'experimental_setup'.
 - Prerequisite: The 'experimental_setup' contains CLEAR limitations or sources of error.
 - 'YES': The answer presents its conclusion as robust WITHOUT mentioning the known limitations.
 - 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
 - 'N/A': The prerequisite is NOT met.

- 5. Formatting & Mechanical Bugs
 - 'flag_incorrect_answer_reference':
 - Evidence Source: 'question', 'suggested_answer', 'ground_truth_answer'.
 - Definition: Whether the provided justification or reasoning identifies the correct answer option(s), BUT then a different option letter is given as the final answer.
 - Prerequisite: The 'question' IS a multiple-choice question (MCQ).

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

- 'YES': The justification or reasoning provided refers to one option letter while discussing the content of another.
- 'NO': The prerequisite IS met, but the conditions for 'YES' are NOT satisfied.
- 'N/A': The prerequisite is NOT met.

Output: Provide your evaluation in the specified JSON format, including the single 'primary_error_category' and the choice ('YES', 'NO', or 'N/A') for every 'flag_'. Note that for some flags the ONLY possible choices as 'YES' and 'NO' (NOT 'N/A'). For each flag, include a brief but clear justification (1-2 sentences) explaining your provided judgment.

[USER]

Fields: outcome_prediction_question, pq_format, experimental_setup, measurement_taken, required_background_knowledge, answer, reasoning_for_answer, clean_gta

Given the following 'experimental_setup' and 'measurements_taken' and 'background_knowledge' (if any):

- 'experimental_setup':
""
{experimental_setup}

- 'measurements_taken':
""
{measurement_taken}

- 'background_knowledge':
""
{required_background_knowledge}

And for the following 'question' (of type 'question_type') and its 'ground_truth_answer':

- 'question_type': {pq_format}

- 'question' (along with choices if applicable):
""
{outcome_prediction_question}

- 'ground_truth_answer':
""
{clean_gta}

Evaluate the following 'suggested_answer' with respect to the provided materials as instructed:

- 'suggested_answer':
""
{answer}

REASONING: {reasoning_for_answer}

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

Prompt for generating responses with background knowledge

[SYS]

Fields: domain, field, experimental_setup, measurement_taken, required_background_knowledge

Instructions: You are tasked with predicting the outcome of a scientific experiment in {domain} and the field of {field} given the provided 'experimental_setup' and 'measurements_taken'. You must analyze the user's scientific 'question' very carefully, and forecast the results AS ACCURATELY AS POSSIBLE given the inputs provided. Each question will have a type (multiple-choice, free-form, numerical) that you must consider when formulating your predictions. Ensure that your predictions are well-reasoned and based on the data provided.

Inputs :

- 'domain': {domain}
- 'field': {field}
- 'experimental_setup': {experimental_setup}
- 'measurements_taken': {measurement_taken}
- 'required_background_knowledge': {required_background_knowledge}

Question Types:

- Multiple-Choice: Choose the most likely outcome from the list of provided options.
- Free-Form: Provide a comprehensive but concise explanation of the expected results.
- Numerical: Predict a specific numerical value of the outcome based on the provided data.

Output: Depending on the 'question_type' provided by the user and based on the provided background knowledge, output the appropriate prediction in the following output fields:

- 'answer'
 - Multiple-Choice: Write ONLY the letter(s) corresponding to the most likely outcome in the 'answer' field (e.g., "X"). If choosing multiple letters (items) is allowed by the 'question' and desired, separate them with commas (e.g., "X, Y, Z").
 - Free-Form: Provide a comprehensive but concise explanation of the expected results.
 - Numerical: Write ONLY the predicted numerical value in the 'answer' field (e.g., "1.234").
- 'reasoning_for_answer': A detailed explanation of how you arrived at your prediction, including any relevant calculations, assumptions, or scientific principles applied.
- 'confidence': Choose between the levels provided. "Confidence" refers to how certain you are about the accuracy of your prediction based on the information provided.
- 'difficulty': Choose between the levels provided. "Difficulty" refers to the complexity of accurately predicting the outcome of the experiment based on the information provided.
- 'feasibility': Choose between the levels provided. "Feasibility" refers to the practicality of predicting the outcome of the experiment WITHOUT conducting it, based on the information provided.
- 'reasoning_for_feasibility': A detailed explanation of how you arrived at your feasibility assessment, considering factors such as experimental design, measurement accuracy, and potential sources of error.

Ensure that your predictions are clear, concise, and directly address the user's scientific 'question'.

[USER]

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

Fields: pq_format, outcome_prediction_question

Answer the following 'question' as accurately as possible:

- 'question_type': {pq_format}
- 'question': {outcome_prediction_question}

Prompt for generating responses without background knowledge

[SYS]

Fields: domain, field, experimental_setup, measurement_taken

Instructions: You are tasked with predicting the outcome of a scientific experiment in { domain} and the field of {field} given the provided 'experimental_setup' and ' measurements_taken'. You must analyze the user's scientific 'question' very carefully, and forecast the results AS ACCURATELY AS POSSIBLE given the inputs provided. Each question will have a type (multiple-choice, free-form, numerical) that you must consider when formulating your predictions. Ensure that your predictions are well-reasoned and based on the data provided.

Inputs :

- 'domain': {domain}
- 'field': {field}
- 'experimental_setup': {experimental_setup}
- 'measurements_taken': {measurement_taken}

Question Types:

- Multiple-Choice: Choose the most likely outcome from the list of provided options.
- Free-Form: Provide a comprehensive but concise explanation of the expected results.
- Numerical: Predict a specific numerical value of the outcome based on the provided data.

Output: Depending on the 'question_type' provided by the user, output the appropriate prediction in the following output fields:

- 'answer'
 - Multiple-Choice: Write ONLY the letter(s) corresponding to the most likely outcome in the 'answer' field (e.g., "X"). If choosing multiple letters (items) is allowed by the 'question' and desired, separate them with commas (e.g., "X, Y, Z").
 - Free-Form: Provide a comprehensive but concise explanation of the expected results.
 - Numerical: Write ONLY the predicted numerical value in the 'answer' field (e.g., "1.234").
- 'reasoning_for_answer': A detailed explanation of how you arrived at your prediction, including any relevant calculations, assumptions, or scientific principles applied.
- 'confidence': Choose between the levels provided. "Confidence" refers to how certain you are about the accuracy of your prediction based on the information provided.
- 'difficulty': Choose between the levels provided. "Difficulty" refers to the complexity of accurately predicting the outcome of the experiment based on the information provided.
- 'feasibility': Choose between the levels provided. "Feasibility" refers to the practicality of predicting the outcome of the experiment WITHOUT conducting it, based on the information provided.
- 'reasoning_for_feasibility': A detailed explanation of how you arrived at your feasibility assessment, considering factors such as experimental design, measurement accuracy, and potential sources of error.

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375

Ensure that your predictions are clear, concise, and directly address the user's scientific 'question'.

[USER]

Fields: pq_format, outcome_prediction_question

Answer the following 'question' as accurately as possible:

- 'question_type': {pq_format}
- 'question': {outcome_prediction_question}

Prompt used for judge

[SYS]

Fields: domain, field, rubric_criteria_lines

Instructions: You are acting as an impartial judge evaluating a suggested answer ('suggested_answer') to a scientific prediction question in the {domain} domain and the field of {field}. Your goal is to determine how well the 'suggested_answer' aligns with the 'ground_truth_answer' based on a set of specific 'rubric_criteria' (a list of ≥ 1 criterion items). Each criterion will need to be evaluated independently. Your evaluation must be objective, rigorous, and strictly based on the provided information. The 'question' was asked given the context information of a scientific experiment as defined by the provided 'experimental_setup' and 'measurements_taken'.

Evaluation Requirements:

1. First, carefully read and understand the scientific context (domain, field) and the specific 'question'. Use the provided 'experimental_setup' and 'measurements_taken' to inform your understanding.
2. Compare the 'suggested_answer' with the 'ground_truth_answer' and reason about the overall correctness and completeness of the 'suggested_answer'.
3. For EACH criterion (INDEPENDENTLY) provided in the 'rubric_criteria' list (could be 1 or more criterion items), you must meticulously assess if the 'suggested_answer' satisfies it ("true" or "false"). The ground truth answer should be used as the reference as the overall correct answer to the 'question'. Provide the output in the corresponding '_satisfied' fields.
4. Your judgment must be objective. Do not introduce external knowledge or make assumptions beyond the provided text.
5. Provide a concise yet clear justification for EACH criterion's determined satisfaction status ("true"/"false") in the corresponding '_reasoning' field.

Inputs:

- 'domain': {domain}
- 'field': {field}
- 'rubric_criteria': Provided below as a list.

Evaluation Criteria:

{rubric_criteria_lines}

Output Format:

You MUST provide your evaluation in a strict JSON format. For each criterion, you will output two fields: one boolean ('_satisfied') and one string ('_reasoning').

[USER]

Fields: outcome_prediction_question, predicted_answer, clean_gta, experimental_setup, measurement_taken

Given the following 'experimental_setup' and 'measurements_taken':

- 'experimental_setup':

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

```

"""
{experimental_setup}
"""

- 'measurements_taken':
  """
  {measurement_taken}
  """

Evaluate the following 'question' with respect to the provided 'suggested_answer' and '
ground_truth_answer' as instructed:
- 'question': {outcome_prediction_question}
- 'suggested_answer': {predicted_answer}
- 'ground_truth_answer': {clean_gta}

```

Prompt to generate responses for questions on background knowledge

```

[SYS]
Fields: domain, field

Instructions: You are tasked with answering questions about a scientific knowledge/facts
in the {domain} domain and the field of {field}. You will be provided with the
experimental setup ('experimental_setup') and the measurements taken ('
measurement_taken') as additional context that are relevant to the questions.
Using this information, you must answer the provided question ACCURATELY
and COMPLETELY.

Output: Provide your accurate and complete answer to each provided question clearly
and concisely. Provide your reasoning for the provided answers in the
corresponding output fields.

[USER]
Fields: bkg_to_qa, experimental_setup, measurement_taken

Given the following 'experimental_setup' and 'measurements_taken':
- 'experimental_setup':
  """
  {experimental_setup}
  """

- 'measurements_taken':
  """
  {measurement_taken}
  """

Answer each of the following questions (each question has a unique hash identifier):
{bkg_to_qa}

```

Prompt to generate questions on background knowledge

```

[SYS]
Fields: domain, field

You are tasked with converting a list of scientific knowledge/fact items in the {domain}
domain and the field of {field} into a set of clear, answerable questions. You will be
provided with the description of the experimental setup and the measurements

```

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

taken, the purpose of the given scientific knowledge/fact items is to help predict the outcome of the experiment. You must create EXACTLY ONE question where the original knowledge/fact is the complete and direct answer. DO NOT MAKE any direct references to the experimental setup, and the measurements taken in the questions.

Output the list of questions and the corresponding original facts in the required JSON format.

[USER]

Fields: experimental_setup, measurement_taken, required_background_knowledge_hashed

Given the following 'experimental_setup' and 'measurements_taken':

- 'experimental_setup':

```
"""
    {experimental_setup}
    """
```

- 'measurements_taken':

```
"""
    {measurement_taken}
    """
```

List of knowledge/fact items to convert:

```
{required_background_knowledge_hashed}
```

Prompt for generating synthetic background knowledge

[SYS]

Fields: domain, field

Instructions: You are tasked with generating relevant background knowledge required for predicting the outcome of the provided scientific experiment in the {domain} domain and the field of {field}. Based on the provided domain, field, experimental setup, and measurements, identify and list 3-6 key scientific principles, facts, or concepts that are essential for predicting the outcome.

Output: Your output must match the required JSON format. Output ONLY a single background knowledge item as an element of the output list (multiple items in the list collectively resulting in multiple pieces of background knowledge). Do NOT output ANY additional comments or text outside in addition to the actual pieces of background knowledge.

Example Output:

```
{
  "generate_bkg": [
    "Background sentence 1.",
    "Background sentence 2."
  ]
}
```

[USER]

Fields: domain, field, experimental_setup, measurement_taken

Please generate the background knowledge for the following experimental direction:

```
- Domain: {domain}
- Field: {field}
- Experimental Setup: {experimental_setup}
- Measurements Taken: {measurement_taken}
```

2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537

Prompt for judging answers to questions on background knowledge

[SYS]

Fields: domain, field

Instructions: You are acting as an impartial judge evaluating a list of answers ('answers') to questions and if those answers capture the corresponding ground truth facts ('ground_truth_facts') for that question in the context of a scientific experiment in the {domain} and the field of {field}. You will also be provided with the experimental setup ('experimental_setup') and measurements taken ('measurements_taken') as additional context that are relevant to the questions. Your goal is to determine if each answer is factually correct and complete (using a coverage metric) based on the provided ground truth facts.

Output: Output your evaluation in the provided JSON format. Each corresponding answer/fact pair is guaranteed to match with a unique hash identifier. For completeness coverage, output a number strictly in the range [0, 1] representing the fraction of ground truth facts that are covered by the answer. For correctness, output "true" if the answer is factually correct with respect to the ground truth facts, and "false" otherwise. Provide a concise yet clear justification for each judgment in the corresponding 'reasoning' fields.

[USER]

Fields: answer_bkg_qa, experimental_setup, measurement_taken, required_background_knowledge_hashed

Given the following 'experimental_setup' and 'measurements_taken':

- 'experimental_setup':

```
"""
    {experimental_setup}
    """
```

- 'measurements_taken':

```
"""
    {measurement_taken}
    """
```

And the following 'ground_truth_facts' (IDs provided in the start of the lines):
{required_background_knowledge_hashed}

Provide your judgments strictly matching the above criteria on the correctness and completeness coverage of each ANSWER against the ground truth (ANSWERS need to be evaluated NOT the ground truth facts):

```
{answer_bkg_qa}
```

Prompt for converting MCQ to FF

[SYS]

Fields: domain, field

Instructions: You are an task with converting multiple-choice questions (MCQ) provided in the {domain} domain and the field of {field} to a free-form question format. You will be provided with the original questions, the multiple-choice options, and the correct answer(s) (potentially multiple), as well as the experimental setup and the measurements taken for the experiment.

Output: Provide the corresponding free-form output question and provide a clear but concise reasoning for the choice and writing of the question. The question must NOT include ANY part from the final MCQ answer and must also not be dependent on the experimental setup or measurements as much as possible. The

2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591

goal is to have a responder answer the output free-form question, and for a judge to then be able to check whether the free-form question was answered correctly and completely or not based on the original correct answer(s) to the original MCQ question. You should also provide an explanation of how a judge would then be able to verify the correctness AND completeness of an answer to the output free-form question given ONLY the original MCQ question and correct answer(s) as well as experimental setup and measurements taken. Questions MUST be clear in scope (not too broad or too narrow), unambiguous, targeted, and end with a question mark.

[USER]
Fields: outcome_prediction_question, experimental_setup, measurement_taken, clean_gta

Given the following 'experimental_setup' and 'measurements_taken':
- 'experimental_setup':
 """
 {experimental_setup}
 """
- 'measurements_taken':
 """
 {measurement_taken}
 """

Convert the following multiple-choice question into a free-form question based on the provided instructions.
{outcome_prediction_question}

Correct answer(s) for this question (NOT to be included in the output free-form question):
{clean_gta}

Provide your output in the specified JSON format, including the new free-form question, your reasoning for constructed it that way, and the explanation for how a judge would verify the correctness and completeness of an answer to the free-form question.