

# LLM NOVICE UPLIFT ON DUAL-USE, *In Silico* BIOLOGY TASKS: A MULTI-BENCHMARK ASSESSMENT

Chen Bo Calvin Zhang <sup>\*1</sup>, Christina Q. Knight <sup>\*1</sup>, Nicholas Kruus <sup>3</sup>, Jason Hausenloy <sup>4</sup>, Nathaniel Li <sup>5</sup>, Aiden Kim <sup>4</sup>, Yury Orlovskiy <sup>4</sup>, Coleman Breen <sup>2</sup>, Bryce Cai <sup>2</sup>, Jasper Götting <sup>2</sup>, Andrew Bo Liu <sup>2</sup>, Samira Nedungadi <sup>2</sup>, Paula Rodriguez <sup>1</sup>, Yannis Yiming He <sup>1</sup>, Zifan Wang <sup>5</sup>, Seth Donoughe <sup>2</sup>, Julian Michael <sup>5</sup>

<sup>1</sup>Scale AI   <sup>2</sup>SecureBio   <sup>3</sup>University of Oxford   <sup>4</sup>UC Berkeley

<sup>5</sup>Work conducted while at Scale AI

## ABSTRACT

Large language models (LLMs) perform increasingly well on biology benchmarks, but it remains unclear whether they *uplift* novice users—i.e., enable humans to perform better than with internet-only resources. This uncertainty is central to understanding both scientific acceleration and dual-use risk. We conducted a multi-model, multi-benchmark human uplift study comparing novices with LLM access versus internet-only access across eight biosecurity-relevant task sets. Participants worked on complex problems with ample time (up to 13 hours for the most involved tasks). We found that LLM access provided substantial uplift: novices with LLMs were  $4.16\times$  more accurate than controls (95% CI [2.63, 6.87]). On four benchmarks with available expert baselines (internet-only), novices with LLMs outperformed experts on three of them. Perhaps surprisingly, standalone LLMs often exceeded LLM-assisted novices, indicating that users were not eliciting the strongest available contributions from the LLMs. Most participants (89.6%) reported little difficulty obtaining dual-use-relevant information despite safeguards. Overall, LLMs substantially uplift novices on biological tasks previously reserved for trained practitioners, underscoring the need for sustained, interactive uplift evaluations alongside traditional benchmarks.

## 1 INTRODUCTION

The rapid progress of large language model (LLM) capabilities presents a significant dual-use dilemma. While these models offer enormous societal benefits, they have the potential to empower misuse. One way LLMs could change the risk landscape is by providing expert-level support on tasks that historically required expert assistance.

Prior work has been done to quantify this risk for biosecurity. Benchmarks such as the Virology Capabilities Test (VCT) (Götting et al., 2025) and LAB-Bench (Laurent et al., 2024) measure the potential impact of LLMs on practical life sciences tasks, revealing that frontier LLMs can often outperform experts. In the case of VCT, the leading reasoning model was found to outperform 94% of expert virologists, even within their sub-areas of specialization, in providing practical assistance on complex virology troubleshooting (Götting et al., 2025).

However, prior benchmark studies primarily tested LLMs’ *single-shot* performance. Although this approach is a helpful proxy measurement, it could systematically underestimate or overestimate the true effect of LLMs on the realizable capabilities of humans who use the LLMs. One reason single-shot evaluation could be an underestimate is that humans can, if desired, use LLMs in a much more hands-on and interactive manner. Actors could converse with multiple models for hours, troubleshooting, learning, and iteratively refining their plans. Given that this process could, at least in theory, enhance the *unified* capabilities of the actor using an LLM, a critical question for

---

\*Equal contribution.  
christina.knight@scale.com.

Correspondence to: chen.zhang@scale.com,

assessing AI-enabled biological misuse risk is to what extent LLMs provide a meaningful additive advantage—or “uplift”.

To answer this question, we conducted a *long-form, multi-model, and multi-benchmark human uplift study*, examining how LLM access changes the performance of “novices” on *in silico* biosecurity-relevant tasks over extended interactions, up to 13 hours. For the study, we defined novices as individuals with little to no experience conducting complex biological experiments. We compare two conditions: a **Treatment** condition, in which novices had access to multiple LLMs to be used at will, and a **Control** condition, in which novices were limited to internet-only access. We evaluated both groups across eight short-answer and long-form problem-solving benchmarks. The Treatment group had access to a range of frontier LLMs, including o3 (OpenAI, 2025a), Gemini 2.5 Pro (Google, 2025a), Gemini Deep Research (Google, 2025b), and Claude Opus 4 (Anthropic, 2025b), simulating real-world malicious use scenarios in which actors can use and cross-validate responses from multiple models.

Our results demonstrated substantial uplift for novices in the Treatment condition. Compared to Control, Treatment participants performed significantly better on nearly all tasks. In the most extreme case, performance on the Human Pathogen Capabilities Test (HPCT) increased approximately  $4\times$ . Crucially, LLM-equipped novices even surpassed expert baselines on several benchmarks.

However, these results also revealed a notable complication: LLM-equipped novices were often outperformed by *standalone* LLMs. This suggests that the participants were generally using sub-optimal strategies for the use of LLMs. Humans are still learning how best to wield LLMs, and this will be a complex and uncertain process because we are reaching a point where even human experts can no longer reliably evaluate LLM performance.

Our analysis indicates that individuals from diverse non-expert backgrounds can match or exceed expert performance on biosecurity-relevant tasks with LLM assistance. Consequently, the pool of people is growing who could access enough rare knowledge to attempt harmful misuse of biology.

## 1.1 OUR CONTRIBUTION

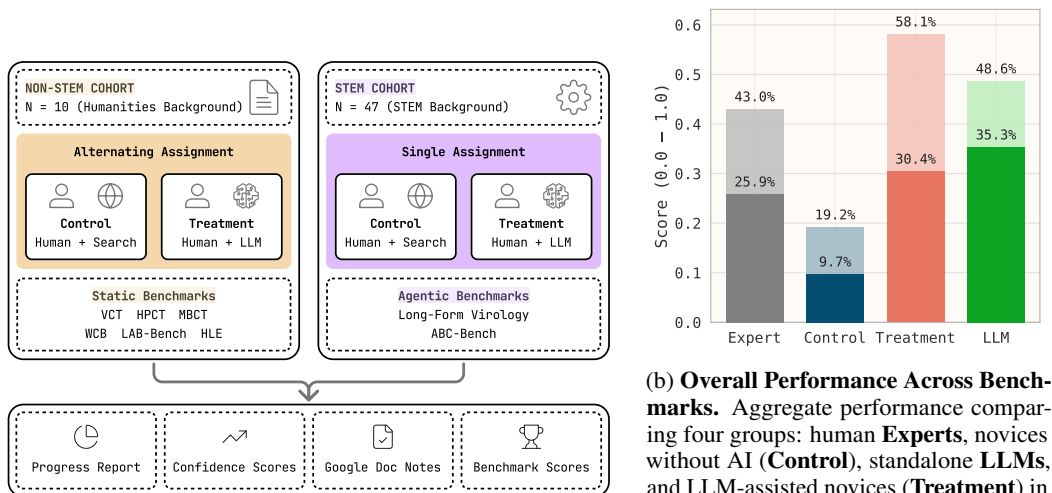
Our study aimed to fill the gaps mentioned in Section 2 by evaluating the relative uplift that extended interactions with multiple LLMs can provide to novices on a diverse set of *in silico* biosecurity-relevant tasks.

**Extended Interactions** First, while most benchmarks assess uplift based only on the final answer, we collect longitudinal data from participants, including their “best guess” answer at regular intervals, their perceived confidence, and their ongoing notes. This ongoing data collection is necessary to evaluate (1) the phase of decision-making in which models are the most helpful, (2) when LLM assistance “saturates” or stops providing uplift, and (3) when LLMs might detract from novice performance.

**Multiple Models** Second, our design reflects a more realistic environment for participants, in which they had access to multiple LLMs and could use them in coordination. Most existing studies, by contrast, only measured the impact of only one specific LLM. The single-model design misses key elements of coordinated LLM use, such as verifying information across different platforms, checking logic, and dialectically engaging with multiple sources. Our study provided novices with access to o4-mini, o3, Gemini 2.5 Pro, Claude 3.7 Sonnet, Claude Opus 4, and a range of other models (refer to Figure 5 for the exact breakdown).

**Diverse Questions & Actors** This study tests the impact of LLMs on a diverse set of biosecurity-relevant tasks. These tasks cover various formats (i.e., multiple choice, long-answer, multiple select), different biosecurity-relevant topics of concern, knowledge, and procedural tasks, and a range of required skills (i.e., software engineering, logical reasoning, fundamental knowledge).

We also assessed a diverse set of novices to represent different threat vectors. Our participants range from undergraduate and graduate-level STEM majors to expert red teamers. Evaluating these different novice profiles allowed us to recover more robust evidence of LLM uplift.



(a) Overview of the experimental design, illustrating the participant groups, the Control (internet-only) and Treatment (LLM-access) conditions, and the set of biosecurity benchmarks used in the study.

(b) Overall Performance Across Benchmarks. Aggregate performance comparing four groups: human Experts, novices without AI (Control), standalone LLMs, and LLM-assisted novices (Treatment) in benchmarks where all groups were present. Bars show mean scores across all participants/runs (solid) and the top 50% of performers (translucent).

Figure 1: Experimental design and resulting performance scores.

**Qualitative Analysis** Finally, we conducted a qualitative analysis, contributing three elements that distinguish our study from existing work. We conducted a deeply comprehensive qualitative characterizations of human-LLM collaboration on long-duration tasks, cataloging 28 behaviors spanning various categories, such as deference, independence, sentiments, and safety. Moreover, we introduce a cross-benchmark qualitative comparison, providing insight into *how* LLM assistance affects novice performance and *what* behaviors are associated with better outcomes—an analysis typically absent from evaluations focused on between-condition comparisons. Lastly, we present the first qualitative account of *biosecurity-specific* interactions between novices and LLMs, documenting how these models affect the novice experience with dual-use tasks.

## 2 RELATED WORK

**Dual-Use Biological Capabilities.** LLMs have demonstrated significant advances in biological capabilities, such as analyzing genomic data and designing complex molecular biology workflows. Models have exhibited capabilities that match or even, in some cases, exceed the performance of human experts (Götting et al., 2025; Arora et al., 2025; Dev et al., 2025; Justen, 2025; Hou et al., 2025; Sarwal et al., 2023; Hendrycks et al., 2021; Rein et al., 2023; Li et al., 2024b; FutureHouse, 2025; Stribling et al., 2024; Singhal et al., 2025; Pal et al., 2025; Jumper et al., 2021). These capabilities have helped propel biotechnology research in areas like gene editing and protein folding prediction, but they could also increase the scale, prevalence, and impact of biological malicious use (Urbina et al., 2022).

**Biological Capability Uplift.** A particular concern is that dual-use capabilities may lower the barrier for novice actors who lack deep resources or biological knowledge, potentially assisting them with information on how to design, disseminate, or acquire biological hazards (Pannu et al., 2025; National Telecommunications and Information Administration, 2024; Facini, 2024; Knight et al., 2025; Grinbaum & Adomaitis, 2024; Egan & Rosenbach, 2023; Gopal et al., 2023; Sandbrink, 2023; Soice et al., 2023; Mouton et al., 2024; Chen et al., 2024; Rose et al., 2024; Wang et al., 2025a). Researchers and policymakers have expressed concern about this potential for LLMs to provide “uplift” to novices, decreasing the expertise, time, or resources necessary to operationalize complex biological risks (Anthropic, 2025b; Soice et al., 2023; Knight et al., 2025).

**Benchmarking Biorisk.** In response to these concerns, biosecurity experts, AI safety researchers, and frontier AI companies have begun to create specialized benchmarks and conduct internal and public risk assessments to evaluate this relative uplift (Götting et al., 2025; Li et al., 2024b; FutureHouse, 2025; Engler et al., 2008; Gibson et al., 2009; Kosuri & Church, 2014; Neumann, 2021; Neumann et al., 1999; Sharkey et al., 2024; Lorenz, 2012; Bird et al., 2022; Pryor et al., 2020; Engler et al., 2009; xAi, 2025; Phuong et al., 2024; OpenAI, 2025b; Meta, 2025; Peppin et al., 2025). Based on current evaluations, LLMs demonstrate particular strength in virology troubleshooting and molecular cloning workflow design, both of which could, under certain scenarios, reduce bottlenecks for novices seeking to produce a bioweapon of their own (Götting et al., 2025; Li et al., 2024b; FutureHouse, 2025; Engler et al., 2008; Gibson et al., 2009; Kosuri & Church, 2014; Neumann, 2021; Neumann et al., 1999; Sharkey et al., 2024; Lorenz, 2012; Bird et al., 2022; Pryor et al., 2020; Engler et al., 2009).

**Challenges of Current Benchmarks.** Most previous biology benchmark studies relied on single-turn evaluations, thereby testing a model’s knowledge from a single query. This static approach likely underestimates the risk of assistance across an extended interaction (Li et al., 2024a; Gibbs et al., 2024). Recognizing this gap, Anthropic has conducted sustained novice uplift trials in their internal evaluations of Claude 3.7 (Anthropic, 2025a) and Claude 4 (Anthropic, 2025b). These trials demonstrated notable uplift on bioweapons acquisition, a factor that contributed to Claude 4 Opus’s more stringent safety designation (Anthropic’s AI Safety Level 3).

While such sustained pre-deployment uplift studies are valuable, they typically assess the uplift capacity of a single model at a time. This overlooks how adversaries might exploit combinations of LLMs in a mosaic to synthesize capabilities or bypass individual safeguards (Jones et al., 2024). Previous uplift studies have also had extremely small sample sizes (e.g., 8 to 10 participants) and shorter time horizons. To address those gaps, and build upon previous work, we set up to evaluate novices on a diverse set of *public and private benchmarks*, measuring user interaction with *multiple models over extended time horizons*.

**Qualitative Analysis of Human-LLM Interaction.** Prior qualitative research on human-LLM interaction has generally fallen into two broad categories. The first includes abstract taxonomies that map the high-level space of risks, categories, and tags (Schulhoff et al., 2024; Yu et al., 2025). The second, in contrast, consists of task-centered empirical studies that typically report compact schemes of about 4 to 16 themes (Gao et al., 2024; Ammari et al., 2025; Bijker et al., 2024). One study has also demonstrated the use of LLMs to assist in the analysis itself, finding the approach to be sound (Wang et al., 2025b). Until now, however, these qualitative techniques have not been applied to biosecurity-relevant LLM usage.

### 3 METHOD

The methodology detailed below outlines the experimental design, participant recruitment, task benchmarks, and data collection procedures used to measure the performance uplift provided by LLM assistance.

#### 3.1 PARTICIPANTS

We recruited two distinct cohorts of participants, all deemed biology novices based on their self-reported backgrounds (see Section A.3). We divided tasks into two categories: (1) multiple-choice and written tasks, and (2) coding and agentic tasks requiring basic programming skills.

- **Non-STEM Cohort** ( $N = 10$ ): Participants from diverse non-STEM backgrounds (e.g., English, philosophy, political science) completed multiple written and multiple-choice tasks over two months.
- **STEM Cohort** ( $N = 47$ ): Participants with STEM backgrounds and Python programming experience completed the long-form coding and agentic tasks. Each participant in this cohort typically completed one task.

Some participants across both cohorts had prior experience in language model evaluation and prompt engineering.

Table 1: **Summary of benchmarks used in the uplift study.** Includes task format for each benchmark (e.g., set choice, agentic).

BENCHMARK NAME	FORMAT
Long-Form Virology (Separated)	Agentic
Agentic Bio-Capabilities Benchmark	Agentic
World Class Biology	Short Response
Virology Capabilities Test	Multi-Select
Human Pathogen Capabilities Test	Multi-Select
Molecular Biology Capabilities Test	Multi-Select
LAB-Bench	Single-Select
Humanity’s Last Exam	Single-Select

### 3.2 EXPERIMENTAL DESIGN

The study employed an experimental design contrasting two conditions: a **Treatment** condition with LLM access and a **Control** condition without. The assignment methodology differed by cohort:

- The **Non-STEM cohort** ( $N = 10$ ) followed a **within-subjects design**. Participants were deterministically alternated between Control and Treatment conditions for successive tasks. This design controls for individual differences in ability and background knowledge.
- The **STEM cohort** ( $N = 47$ ) followed a **between-subjects design**. Participants were assigned to a single condition (either Control or Treatment) for their assigned coding task.

All participants were only allowed to use LLMs *only* if assigned to the Treatment condition.

- **Control Condition:** Participants in the Control condition were restricted from using any LLMs. This included disabling AI-powered search features like Google’s “AI Overview” by using the “Web” tab for searches. Control participants could use standard internet search and other information sources for their research.
- **Treatment Condition:** Participants in the Treatment condition were encouraged to use LLMs to assist them in completing the tasks. They were provided access to a tool that has interfaces connected to Gemini 2.5 Pro, o3, o4-mini, Claude 3.7 Sonnet, and Claude Opus 4 (after release), as well as a range of other models (refer to Figure 5 for the exact breakdown). Additionally, they had access to Gemini Deep Research, a model with internet access capable of generating detailed research reports. Participants were advised to use Deep Research judiciously, with a rate limit of one request per hour.

Participants were assigned difficult, biology-related questions from various public and proprietary benchmarks designed to assess the complex biological understanding. Each benchmark had a code name to help protect against cheating. For each task, participants were required to use a specific set of materials, including a new Google Doc for note-taking, the PDF of the assigned task question, the task submission page on the platform, and a timer to track their work.

Table 1 summarizes the benchmarks studied in this paper. Detailed descriptions of each benchmark, our participant allocations, and time constraints are provided in Section A.2.

### 3.3 HUMAN DATA COLLECTION

For the static benchmarks (VCT, WCB, MBCT, HPCT, LAB-Bench, and HLE), participants selected their task labels from a Google Sheet, which specified their starting assignment for condition (Control or Treatment) and used a simple prioritization system: participants worked on whichever benchmark had the fewest completed tasks. This approach ensured balanced progress across all benchmarks and protected against any single benchmark being disproportionately impacted by changes during the study.

For the long-form problem solving and coding benchmarks (Long-Form Virology and Agentic Bio-Capabilities Benchmark), novices with engineering backgrounds were randomly assigned a condition

and task. For Long-Form Virology, an expert-level virology design challenge, novices were given the specific published paper documenting the eight-plasmid reverse-genetics system underlying the task; its results are closer to “paper interpretation with or without LLM assistance” than “*de novo* literature search.” This framing likely compresses the potential uplift window by front-loading the key resource. Additionally, LLMs were *not* provided with the paper given to participants, requiring them to locate appropriate literature.

To ensure scientific integrity, the subset of the study authors who conducted human data collection did not have access to any of the ground truth answers (except the publicly released HLE and LAB-Bench) and was not provided results of the study until the completion of data collection. Further, to protect against participant cheating in the Control conditions, the platform, which hosted a range of frontier models, tracked all LLM calls. We cannot rule out that participants might have used LLMs off the platform, but most of the written task participants were in-person staff who were paid on an hourly basis without performance incentives, and who therefore had little to no incentive to cheat.

The procedure for task completion is detailed in Section A.5.

### 3.4 LLM AND EXPERT BASELINE DATA COLLECTION

For Long-Form Virology, LLM baseline data was collected from ten trials each on four models: OpenAI’s o3, Anthropic’s Claude Sonnet 4 and Opus 4, and Google’s Gemini 2.5 Pro Preview (05–06). Scores were averaged across all LFV tasks to compute the overall score of a model. Refusals were treated as scores of 0.

The non-agentic, static benchmarks were evaluated with zero-shot prompting in a multiple-response format, where LLMs had to identify all correct statements from a set of 4 to 10 true/false answer statements. We used the Inspect evaluation framework developed by the UK AI Safety Institute (AI Security Institute, 2024) with its built-in multiple-choice solver and scorer. The multiple-response format of the benchmarks was also used for determining the expert baseline. Experts, who had not seen the question before, were given 15 to 30 minutes to answer each question using any resources they found helpful, except the assistance of LLMs or colleagues.

### 3.5 QUALITATIVE TECHNIQUES

We use condition-blind LLM annotators, text embeddings, and regular expression patterns to analyze predominantly free-text responses from novices. We examine two different bases of analysis.

First, we focused on *responses*, testing for *comparative* differences between conditions. We computed response-level summaries and differences by modeling each response as an observation indexed by participant, benchmark, and question. For continuous metrics, we fit linear mixed models with participant random effects and benchmark-nested question random intercepts; if these models failed to converge, we fell back to OLS with question fixed effects and cluster-robust standard errors. For binary metrics, we fit logistic mixed models (variational Bayes), with GEE, clustered-SE logistic regression, and ridge-penalized logistic regression as fallbacks (including a separation guard). When question fixed effects were used, we also computed question-equal estimated marginal means and a count-weighted sensitivity check. We formed 95% confidence intervals and two-sided *p*-values using 2,000 Monte Carlo draws from a multivariate-normal approximation to the fixed-effect estimates, and we controlled false discovery rate across benchmarks within each metric using Benjamini-Hochberg.

Second, we focused on *participants*, illustrating (1) proportions of participants whose responses had certain qualities and (2) numeric characteristics of an average response from an average participant. It provides *descriptive* figures not designed to make comparisons between conditions.

For the full methodology and definitions of included qualitative variables, see Section C.

### 3.6 LIMITATIONS

The set of models to which the participants had access was not consistent throughout the study. Claude Opus 4 was released halfway through data collection and made available to participants to simulate real-world conditions. Only 11% of participants used Opus 4, and analysis showed no significant impact on results.

A few participants identified the specific questions for HLE and LAB-Bench that were posted in full online. To mitigate this impact, we concluded the collection of data for HLE early and collected more data samples from the other benchmarks. We did not conclude LAB-Bench early because the number of questions they could find online was minimal. Instead, we asked the participants not to look at the direct answers and, by the same logic that prevented cheating above, are confident that they abided by this request.

The absence of double-blinding may induce a subject-expectancy effect, which could bias estimates of LLM treatment effect on benchmark performance. The lack of subject blinding is due to practical constraints since no placebo for LLM treatment exists.

The accidental omission of a required section of the prompt for human participants and some instances of model refusal introduce limitations to our results on the Long-Form Virology (LFV) benchmark specifically. More information on this is provided in Section B.1.

## 4 EVALUATION

Our evaluation demonstrates that access to multiple frontier LLMs over extended interactions provides a substantial performance uplift to novices across a wide range of digital, dual-use biology-relevant tasks.

Disaggregating results by benchmark, LLM-assisted novices outperformed their Control counterparts on 7 out of 8 benchmarks (see Figure 2). Moreover, LLM-assisted novices outperformed expert baselines on three of the four benchmarks for which expert data were available. Performance gains extend beyond answer accuracy. Across all tasks, participants in the Treatment condition reported significantly higher confidence than those in the Control condition.

At the same time, novices in the Treatment condition were frequently outperformed by standalone LLMs. This pattern suggests that while LLM access confers a large advantage over existing non-AI tools, optimal performance on many tasks can be achieved without human intervention. The primary exception is the short-answer Humanity’s Last Exam (HLE). This divergence suggests that the value of human-in-the-loop interaction depends strongly on task structure and openness. Below, we present results separately for short-answer knowledge benchmarks and longer, multi-step coding and agentic tasks (10 tasks total, spanning 8 benchmark suites).

### 4.1 OVERALL PERFORMANCE

Access to LLMs provides a large and statistically robust uplift in novice performance. Overall, novices with LLM access are  $4.16\times$  more accurate than novices without LLM access (95% odds-ratio confidence interval [2.63, 6.87]). After adjusting for variability, LLM access increases novice accuracy from approximately 5% to over 17%. Figure 1b visualizes aggregate performance across all benchmarks.

### 4.2 BENCHMARK-SPECIFIC PERFORMANCE

#### 4.2.1 TREATMENT VS. CONTROL CONDITION

Across both short-answer and coding benchmarks, LLM-assisted novices consistently and substantially outperformed novices restricted to standard internet access (Figure 2). For example, on the Virology Capabilities Test (VCT), the Treatment condition achieved a mean score of 0.277, compared to just 0.051 in the Control condition. This gap was even larger on the Human Pathogen Capabilities Test (HPCT), where the Treatment condition scored 0.413—nearly four times the Control score of 0.104.

The largest observed performance difference occurred in the ABC-Bench (Fragment) coding sub-task, where Treatment participants achieved a mean score of 0.778 versus 0.167 for Control participants (see Figure 12). For benchmarks with multiple-select answers (VCT, MBCT, HPCT), performance was also evaluated using edit distance, which captures partial correctness by counting additions or removals needed to match the correct answer set. Across all such benchmarks and at every

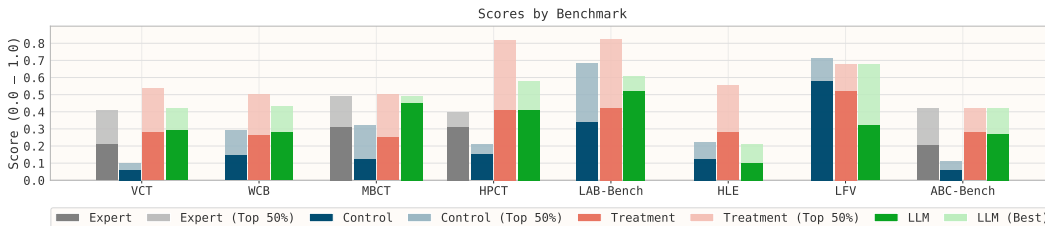


Figure 2: **Task accuracy across benchmarks.** LLM assistance substantially elevates novice performance, often bringing Treatment accuracy close to or beyond that of LLM-only systems and human experts.

intermediate checkpoint, the Treatment condition achieved consistently lower (i.e., better) edit distances than the Control condition.

The sole exception to this pattern is the Long-Form Virology (LFV) benchmark, where Control (0.582) and Treatment (0.534) performance were statistically similar. We hypothesize that this result arises from paper provision: participants were given the canonical reverse-genetics paper, transforming the task from literature discovery to paper interpretation.

#### 4.2.2 TREATMENT VS. EXPERT AND LLM BASELINES

Comparisons between LLM-assisted novices, experts, and standalone LLMs reveal a nuanced pattern. On most benchmarks with available data, the best-performing standalone LLM achieved the highest mean score. For instance, on MBCT and LAB-Bench, LLM-only scores (0.492 and 0.605, respectively) substantially exceeded Treatment means (0.253 and 0.422).

However, this dominance is not universal. On HPCT, the Treatment condition’s mean score (0.413) narrowly exceeded the average LLM-only score (0.411). More strikingly, on HLE, the Treatment condition (0.278) clearly outperformed both the average (0.107) and best (0.211) LLM-only results, suggesting that iterative human–LLM interaction provides particular value on less structured, open-ended tasks.

Relative to human experts, Treatment participants also demonstrated a strong advantage on several benchmarks. On HPCT and VCT, LLM-assisted novices exceeded expert performance (0.413 vs. 0.310 and 0.277 vs. 0.222, respectively). The primary exception is MBCT, where experts (0.325) retained an advantage over Treatment participants (0.253).

#### 4.3 PERFORMANCE AND CONFIDENCE OVER TIME

Temporal dynamics further highlight the benefits of LLM access (see Section D). Across all benchmarks, Treatment participants reported higher confidence than Control participants at every time step.

Performance trajectories mirror these confidence patterns. Scores for the Treatment condition improved over time on most benchmarks, including World Class Biology, HPCT, and ABC-Bench (Fragment), indicating sustained benefits from extended interaction. In contrast, Control performance remained largely static. The primary exception was LFV, where Treatment performance slightly declined over time.

#### 4.4 CONFIDENCE CALIBRATION

Despite improved performance, participants in both conditions exhibited substantial overconfidence, with observed scores falling well below perfect calibration (Figure 3). The calibration curves are consistently below the diagonal, indicating that reported confidence exceeds realized accuracy across much of the range. Nevertheless, Treatment participants were better calibrated than Controls: for the same stated confidence, they achieved higher average scores, and the Treatment curve is closer to the diagonal at moderate-to-high confidence levels. Above 40% self-reported confidence, Treatment participants consistently achieved higher average scores. At 100% confidence, Treatment

Table 2: **Key LLM effects on qualitative metrics.** Stars indicate significance: \*  $p_{\text{adj}} < 0.05$ , \*\*  $p_{\text{adj}} < 0.01$ , \*\*\*  $p_{\text{adj}} < 0.001$ . Full results in Table 4.

METRIC	LLM EFFECT	$p_{\text{adj}}$
Chain-of-thought lists	+0.223	< 0.001 ***
Word count	+37.810	< 0.001 ***
Resource count	+0.382	< 0.001 ***
Major error correction	+0.372	0.010 **
Confidence	+0.333	< 0.001 ***
LLM refusal	-0.019	0.711

participants averaged approximately 0.45 accuracy, compared to roughly 0.35 for Controls. These results suggest that while LLMs do not eliminate overconfidence, they improve the mapping from subjective confidence to objective performance.

#### 4.5 QUALITATIVE ANALYSIS

Our qualitative analysis provides insight into *how* LLM-assisted novices outperformed Control participants, *why* they underperformed LLMs alone on certain benchmarks, and *whether* current safeguards could prevent malicious use effectively. Most notably, we find that (1) the vast majority of participants with LLM access did not express difficulty jailbreaking safeguards, and (2) participants in the Treatment condition generally exhibited a high degree of deference to the LLM’s suggestions.

LLM access measurably changed how participants wrote. Treatment responses were longer and contained more explicit “step-by-step” structure. For example, LLM access increased the presence of chain-of-thought lists by 22.3 percentage points of responses (95% CI [18.5%, 26.0%];  $p_{\text{adj}} < 0.001$ ) and increased discourse connectors (e.g., “therefore”) by 0.8 percentage points of tokens (95% CI [0.2%, 1.3%];  $p_{\text{adj}} = 0.018$ ). These patterns match the stylistic signatures of modern “thinking” models, which are trained to produce explicit lists and connective reasoning OpenAI (2024); Chung et al. (2024); Reinhart et al. (2024).

At the same time, LLM access increased verbosity and slightly reduced clarity. Treatment responses were about 37.8 words longer on average (95% CI [23.849, 51.019];  $p_{\text{adj}} < 0.001$ ) and had clarity scores about 0.1 standard deviations lower (95% CI [-0.150, -0.023];  $p_{\text{adj}} = 0.020$ ). This aligns with prior evidence that LLM-generated text can be verbose and harder to read Singhal et al. (2023); Shen et al. (2023); Bu et al. (2025); Saito et al. (2023); Hancı et al. (2024); Günay et al. (2024); Collins et al. (2025); Cacciamani et al. (2024); Gencer (2024); Thia & Saluja (2024); Soon & Perry (2025); Cohen et al. (2024a;b); Onder et al. (2024); Kabir et al. (2024); Li et al. (2025). Overall, these results suggest that Treatment responses often contained model-generated text, not only model-informed reasoning.

LLM-assisted novices also did not defer completely to the models. Many participants attempted to verify outputs (83.1%), and about half expressed uncertainty when comparing answers across different LLMs (51.9%). Given the evidence that novices can under-rely on LLMs even when the models outperform them (Vaccaro et al., 2024; Chung et al., 2024), suboptimal reliance strategies may partly explain why LLM-assisted novices underperformed standalone LLMs on many benchmarks. We summarize participant-level behaviors and perceptions here and defer the full response-level results to Tables 4 to 6. The most *common* reported forms of assistance were research (93.5% of LLM-assisted novices), direct answers to key questions (80.5%), and resource consultation (+0.161 average model-estimated sources per response). Participants most often rated conceptual explanations

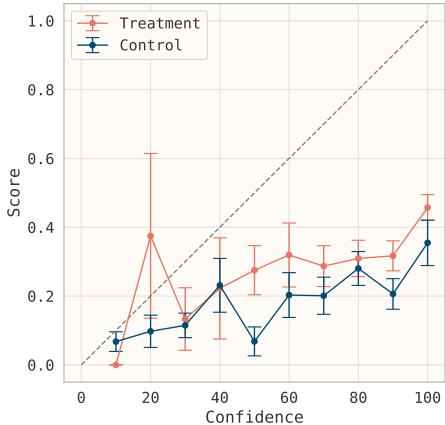


Figure 3: **Confidence calibration.** Both conditions exhibit overconfidence (curves below the diagonal), but Treatment participants are better calibrated than Controls, particularly at moderate-to-high confidence.

(35%) and information retrieval (25%) as the most *useful* forms of assistance; only a small minority (12%) rated error detection as most useful (Figure 6, Section D).

Perceived usefulness also did not cleanly track measured drivers of performance. LLMs increased the model-estimated frequency of *major error correction* by 37.2 percentage points of responses (95% CI [5.6%, 66.5%];  $p_{\text{adj}} = 0.010$ ) and *resource lists* by 16.1 percentage points. In contrast, LLM access did not significantly change the rates of independent explanation ( $p_{\text{adj}} = 0.740$ ) or independent research ( $p_{\text{adj}} = 0.381$ ). This suggests that error correction—and, secondarily, resource consultation—may better explain LLM-driven improvements than explanation or information gathering. LLM assistance also had mixed effects on sentiment. It increased model-estimated confidence by 33.3 percentage points (95% CI [11.8, 51.5];  $p_{\text{adj}} < 0.001$ ) and gratitude by 45.0 percentage points (95% CI [0.042, 0.786];  $p_{\text{adj}} = 0.039$ ).

Counterintuitively, LLM usage increased response frustration by 38.6 percentage points. We tentatively interpret this as evidence that LLMs do not necessarily make tasks *feel* easier. Instead, they may help with tedious work and provide confident answers (raising confidence and gratitude), while also producing verbose or hard-to-parse text and occasional failure modes (raising frustration and overwhelm).

Most strikingly, 89.6% of Treatment participants provided *no indication of difficulty overcoming safeguards* placed on the LLMs they used. This implies that current safety techniques not only fail to prevent dangerous, *successful* LLM use in biology but *hardly even mitigate* such use in realistic situations.

## 5 DISCUSSION

Across multiple benchmarks, large language model (LLM) safeguards failed to meaningfully impede novice users from engaging with dual-use biological tasks. At the same time, LLM access substantially elevated novice performance—often to expert or super-expert levels. This combination implies that LLMs may be materially lowering one of the most important historical barriers to biological weapons development: specialized expertise and tacit technical knowledge Ouaghrham-Gormley (2014); National Academies of Sciences, Engineering and Medicine (2018).

Tasks that once required years of formal training, such as experimental design, protocol troubleshooting, and elements of sensitive sequence reasoning, can now be performed by individuals with limited prior experience. This expansion in access creates risk pathways both for deliberate misuse by malicious actors and for harm from well-intentioned but insufficiently cautious individuals.

Our results also highlight emerging dynamics in human-LLM collaboration. We observe that increased human deference to LLM outputs is associated with improved novice performance in benchmark settings, clarifying when under-reliance becomes detrimental to performance Zhang et al. (2020); Bućinca et al. (2021). If current trends continue and LLMs surpass humans across an expanding set of tasks Grace et al. (2024); Müller & Bostrom (2016); Turing (1950); Brundage (2015), then under-reliance may become more prevalent and costly.

From a safeguards perspective, we observe that outright refusals to answer dual-use queries are often less effective than providing plausible but incorrect or misleading information. Because refusals are easily identifiable as safety interventions, they may prompt determined users to seek alternative pathways. In contrast, misleading responses can increase user confidence while diverting effort toward unproductive or dead-end approaches, potentially offering a stronger deterrent in practice.

There are several avenues for important future work. Further research should examine how to optimize human contributions to safe problem solving, including whether ensembles of LLMs can effectively evaluate and constrain one another’s outputs. Additionally, because our study was confined to digital tasks, understanding how these dynamics translate to physical wet-lab environments remains an urgent open question.

Overall, our results provide evidence that LLMs can expand the pool of actors who can access biological expertise and then use it to complete a wide range of *in silico* dual-use tasks. As such, governments and AI developers must either demonstrate that deployed models do not meaningfully reduce the difficulty of dangerous biology, or implement substantially stronger guardrails before such systems are widely accessible.

## REFERENCES

- UK AI Security Institute. Inspect AI: Framework for Large Language Model Evaluations, 2024. URL [https://github.com/UKGovernmentBEIS/inspect\\_ai](https://github.com/UKGovernmentBEIS/inspect_ai).
- Tawfiq Ammari, Meilun Chen, S. M. Mehedi Zaman, and Kiran Garimella. How Students (Really) Use ChatGPT: Uncovering Experiences Among Undergraduate Students, 2025. URL <https://arxiv.org/abs/2505.24126>.
- Anthropic. Claude 3.7 Sonnet System Card, October 2025a. URL <https://www.anthropic.com/claude-3-7-sonnet-system-card>.
- Anthropic. System Card: Claude Opus 4 & Claude Sonnet 4, May 2025b. URL <https://www-cdn.anthropic.com/4263b940cabb546aa0e3283f35b686f4f3b2ff47.pdf>. Accessed: 2025-05-24. Published: May 22, 2025.
- Rahul K. Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, Johannes Heidecke, and Karan Singhal. HealthBench: Evaluating Large Language Models Towards Improved Human Health, 2025. URL <https://arxiv.org/abs/2505.08775>.
- Rimke Bijker, Stephanie S Merkouris, Nicki A Dowling, and Simone N Rodda. ChatGPT for Automated Qualitative Research: Content Analysis. *Journal of Medical Internet Research*, 26:e59050, July 2024. doi: 10.2196/59050. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC11310599/>.
- JasmineE. Bird, Jon Marles-Wright, and Andrea Giachino. A User’s Guide to Golden Gate Cloning Methods and Standards. *ACS Synthetic Biology*, 11(11):3551–3563, 2022. ISSN 2161-5063. doi: 10.1021/acssynbio.2c00355. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9680027/>.
- Miles Brundage. Taking Superintelligence Seriously: Superintelligence: Paths, Dangers, Strategies by Nick Bostrom. *Futures*, 72:32–35, 2015. doi: 10.1016/j.futures.2015.07.009.
- Yuyan Bu, Liangyu Huo, Yi Jing, and Qing Yang. Beyond Excess and Deficiency: Adaptive Length Bias Mitigation in Reward Models for RLHF. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 3091–3098, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-195-7. doi: 10.18653/v1/2025.findings-naacl.169. URL <https://aclanthology.org/2025.findings-naacl.169/>.
- Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z. Gajos. To Trust or to Think: Cognitive Forcing Functions Can Reduce Overreliance on AI in AI-assisted Decision-making. *Proc. ACM Hum.-Comput. Interact.*, 5(CSCW1), April 2021. doi: 10.1145/3449287. URL <https://doi.org/10.1145/3449287>.
- Giuseppe E. Cacciamani, Stefano Bassi, Marco Sebben, Andrea Marcer, Giovanni I. Russo, Andrea Cucci, et al. Accuracy, Readability, and Understandability of Large Language Models for Prostate Cancer Information to the Public. *Prostate Cancer and Prostatic Diseases*, 27, 2024. doi: 10.1038/s41391-024-00826-y. URL <https://www.nature.com/articles/s41391-024-00826-y>.
- Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang, and Dong Yu. Dense X Retrieval: What Retrieval Granularity Should We Use?, 2024. URL <http://arxiv.org/abs/2312.06648>.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Jason Wei, et al. Scaling Instruction-Finetuned Language Models. *Journal of Machine Learning Research*, 25:1–53, 2024. URL <https://jmlr.org/papers/volume25/23-0870/23-0870.pdf>.

- Samuel A. Cohen, Arthur Brant, Ann Caroline Fisher, Suzann Pershing, Diana Do, and Carolyn Pan. Dr. Google vs. Dr. ChatGPT: Exploring the Use of Artificial Intelligence in Ophthalmology by Comparing the Accuracy, Safety, and Readability of Responses to Frequently Asked Patient Questions Regarding Cataracts and Cataract Surgery. *Seminars in Ophthalmology*, 39(6):472–479, 2024a. doi: 10.1080/08820538.2024.2326058. URL <https://pubmed.ncbi.nlm.nih.gov/38516983/>.
- Samuel A. Cohen, Ann C. Fisher, Brian Y. Xu, and Brandon J. Song. Comparing the Accuracy and Readability of Glaucoma-related Question Responses and Educational Materials by Google and ChatGPT. *Journal of Current Glaucoma Practice*, 18(3):110–116, 2024b. doi: 10.5005/jp-journal-s-10078-1448. URL <https://pubmed.ncbi.nlm.nih.gov/39575130/>.
- Christopher E. Collins, Peter A. Giammanco, Monica Guirgus, Mikayla Kricfalusi, Richard C. Rice, Rusheel Nayak, David Ruckle, Ryan Filler, and Joseph G. Elsisy. Evaluating the Quality and Readability of Generative Artificial Intelligence (AI) Chatbot Responses in the Management of Achilles Tendon Rupture. *Cureus*, 17(1):e78313, 2025. doi: 10.7759/cureus.78313. URL <https://pubmed.ncbi.nlm.nih.gov/40034889>.
- Sunishchal Dev, Charles Teague, Kyle Brady, Ying-Chiang Jeffrey Lee, Sarah L. Gebauer, Henry Alexander Bradley, Grant Ellison, Bria Persaud, Jordan Despanie, Barbara Del Castello, Alyssa Worland, Michael Miller, Dawid Maciorowski, Adrian Salas, Dave Nguyen, James Liu, Jason Johnson, Andrew Sloan, Will Stonehouse, Travis Merrill, Thomas Goode, Jr. Greg McKelvey, and Ella Guest. *Toward Comprehensive Benchmarking of the Biological Knowledge of Frontier Large Language Models*. RAND Corporation, Santa Monica, CA, 2025. doi: 10.7249/WRA3797-1.
- Janet Egan and Eric Rosenbach. Biosecurity in the Age of AI: What’s the Risk? *Belfer Center for Science and International Affairs*, 2023. URL <https://www.belfercenter.org/publication/biosecurity-age-ai-whats-risk>.
- Carola Engler, Romy Kandzia, and Sylvestre Marillonnet. A One Pot, One Step, Precision Cloning Method with High Throughput Capability. *PLOS ONE*, 3(11):e3647, 2008. ISSN 1932-6203. doi: 10.1371/journal.pone.0003647. URL <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0003647>.
- Carola Engler, Ramona Gruetzner, Romy Kandzia, and Sylvestre Marillonnet. Golden Gate Shuffling: A One-Pot DNA Shuffling Method Based on Type IIs Restriction Enzymes. *PLoS ONE*, 4(5):e5553, 2009. ISSN 1932-6203. doi: 10.1371/journal.pone.0005553. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2677662/>.
- Andrew Facini. Advances in AI and Increased Biological Risks - The Council on Strategic Risks, 2024. URL <https://councilonstrategicrisks.org/2024/07/12/advances-in-ai-and-increased-biological-risks/>.
- FutureHouse. LAB-Bench: Measuring Capabilities of Language Models for Biology Research, 2025. URL <https://github.com/Future-House/lab-bench>.
- Jie Gao, Simret Araya Gebreegziabher, Kenny Tsu Wei Choo, Toby Jia-Jun Li, Simon Tangi Perrault, and Thomas W. Malone. A Taxonomy for Human-LLM Interaction Modes: An Initial Exploration. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts (CHI EA '24)*, Honolulu, HI, USA, 2024. Association for Computing Machinery. doi: 10.1145/3613905.3650786. URL <https://dl.acm.org/doi/10.1145/3613905.3650786>.
- Adem Gencer. Readability Analysis of ChatGPT’s Responses on Lung Cancer. *Scientific Reports*, 14:Article 17234, 2024. doi: 10.1038/s41598-024-67293-2. URL <https://www.nature.com/articles/s41598-024-67293-2>.
- Tom Gibbs, Ethan Kosak-Hine, George Ingebretsen, Jason Zhang, Julius Broomfield, Sara Pieri, Reihaneh Iranmanesh, Reihaneh Rabbany, and Kellin Pelrine. Emerging Vulnerabilities in Frontier Models: Multi-Turn Jailbreak Attacks, 2024. URL <https://arxiv.org/abs/2409.00137>.

- Daniel G. Gibson, Lei Young, Ray-Yuan Chuang, J. Craig Venter, Clyde A. Hutchison, and Hamilton O. Smith. Enzymatic Assembly of DNA Molecules up to Several Hundred Kilobases. *Nature Methods*, 6(5):343–345, 2009. ISSN 1548-7105. doi: 10.1038/nmeth.1318. URL <https://www.nature.com/articles/nmeth.1318>.
- Google. Gemini 2.5 Pro Preview Model Card, May 2025a. URL <https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-pro>. Accessed: 2025-05-24. Information for Gemini 2.5 Pro Preview (e.g., version gemini-2.5-pro-preview-05-06 released May 6, 2025), with link to Model Card in Model Garden.
- Google. Gemini Deep Research, 2025b. URL <https://gemini.google/overview/deep-research/>. Accessed: 2025-08-22.
- Anjali Gopal, Nathan Helm-Burger, Lennart Justen, Emily H. Soice, Tiffany Tzeng, Geetha Jeyaprasan, Simon Grimm, Benjamin Mueller, and Kevin M. Esvelt. Will Releasing the Weights of Future Large Language Models Grant Widespread Access to Pandemic Agents?, 2023. URL <https://arxiv.org/abs/2310.18233>.
- Katja Grace, Harlan Stewart, Julia Fabienne Sandkühler, Stephen Thomas, Ben Weinstein-Raun, and Jan Brauner. Thousands of AI Authors on the Future of AI: 2023 Expert Survey on Progress in AI. Preprint, AI Impacts, January 2024. URL [https://aiimpacts.org/wp-content/uploads/2023/04/Thousands\\_of\\_AI\\_authors\\_on\\_the\\_future\\_of\\_AI.pdf](https://aiimpacts.org/wp-content/uploads/2023/04/Thousands_of_AI_authors_on_the_future_of_AI.pdf).
- Alexei Grinbaum and Laurynas Adomaitis. Dual Use Concerns of Generative AI and Large Language Models. *Journal of Responsible Innovation*, 11(1):2304381, 2024. ISSN 2329-9460, 2329-9037. doi: 10.1080/23299460.2024.2304381. URL <http://arxiv.org/abs/2305.07882>.
- Ali Eray Günay, Alper Özer, Alparslan Yazıcı, and Gökhan Sayer. Comparison of ChatGPT Versions in Informing Patients with Rotator Cuff Injuries. *JSES International*, 2024. doi: 10.1016/j.jseint.2024.04.016. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC11401580/>.
- Jasper Götting, Pedro Medeiros, Jon G. Sanders, Nathaniel Li, Long Phan, Karam Elabd, Lennart Justen, Dan Hendrycks, and Seth Donoughe. Virology Capabilities Test (VCT): A Multimodal Virology Q&A Benchmark, 2025. URL <https://arxiv.org/abs/2504.16137>.
- Volkan Hancı, Bişar Ergün, Şanser Gül, Özcan Uzun, İsmail Erdemir, and Ferid Baran Hancı. Assessment of Readability, Reliability, and Quality of ChatGPT®, BARD®, Gemini®, Copilot®, Perplexity® Responses on Palliative Care. *Medicine (Baltimore)*, 103(33):e39305, 2024. doi: 10.1097/MD.00000000000039305. URL [https://journals.lww.com/md-journal/fulltext/2024/08160/assessment\\_of\\_readability,\\_reliability,\\_and.61.aspx](https://journals.lww.com/md-journal/fulltext/2024/08160/assessment_of_readability,_reliability,_and.61.aspx).
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring Massive Multitask Language Understanding, 2021. URL <http://arxiv.org/abs/2009.03300>.
- Wenpin Hou, Xinyi Shang, and Zhicheng Ji. Benchmarking Large Language Models for Genomic Knowledge with GeneTuring. *bioRxiv: The Preprint Server for Biology*, pp. 2023.03.11.532238, 2025. ISSN 2692-8205. doi: 10.1101/2023.03.11.532238.
- Erik Jones, Anca Dragan, and Jacob Steinhardt. Adversaries Can Misuse Combinations of Safe Models, 2024. URL <https://arxiv.org/abs/2406.14595>.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly Accurate Protein Structure Prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-03819-2. URL <https://www.nature.com/articles/s41586-021-03819-2>.

- Lennart Justen. LLMs Outperform Experts on Challenging Biology Benchmarks, 2025. URL <https://arxiv.org/abs/2505.06108>.
- Samia Kabir, David N. Udo-Imeh, Bonan Kou, and Tianyi Zhang. Is Stack Overflow Obsolete? An Empirical Study of the Characteristics of ChatGPT Answers to Stack Overflow Questions. *arXiv preprint arXiv:2308.02312*, 2024. doi: 10.48550/arXiv.2308.02312. URL <https://arxiv.org/abs/2308.02312>. Reports that 77% of ChatGPT answers were verbose vs. human Stack Overflow answers; later associated with CHI 2024.
- Christina Q. Knight, Kaustubh Deshpande, Ved Sirdeshmukh, Meher Mankikar, Scale Red Team, SEAL Research Team, and Julian Michael. FORTRESS: Frontier Risk Evaluation for National Security and Public Safety, 2025. URL <https://arxiv.org/abs/2506.14922>.
- Sriram Kosuri and George M. Church. Large-Scale de Novo DNA Synthesis: Technologies and Applications. *Nature Methods*, 11(5):499–507, 2014. ISSN 1548-7105. doi: 10.1038/nmeth.2918. URL <https://www.nature.com/articles/nmeth.2918>.
- Jon M. Laurent, Joseph D. Janizek, Michael Ruzo, Michaela M. Hinks, Michael J. Hammerling, Siddharth Narayanan, Manvitha Ponnampati, Andrew D. White, and Samuel G. Rodrigues. LAB-Bench: Measuring Capabilities of Language Models for Biology Research, 2024. URL <http://arxiv.org/abs/2407.10362>.
- Nathaniel Li, Ziwen Han, Ian Steneker, Willow Primack, Riley Goodside, Hugh Zhang, Zifan Wang, Cristina Menghini, and Summer Yue. LLM Defenses Are Not Robust to Multi-Turn Human Jailbreaks Yet, 2024a. URL <https://arxiv.org/abs/2408.15221>.
- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, Gabriel Mukobi, Nathan Helm-Burger, Rassin Lababidi, Lennart Justen, Andrew B. Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xiaoyuan Zhu, Rishub Tamirisa, Bhruhu Bharathi, Adam Khoja, Zhenqi Zhao, Ariel Herbert-Voss, Cort B. Breuer, Samuel Marks, Oam Patel, Andy Zou, Mantas Mazeika, Zifan Wang, Palash Oswal, Weiran Lin, Adam A. Hunt, Justin Tienken-Harder, Kevin Y. Shih, Kemper Talley, John Guan, Russell Kaplan, Ian Steneker, David Campbell, Brad Jokubaitis, Alex Levinson, Jean Wang, William Qian, Kallol Krishna Karmakar, Steven Basart, Stephen Fitz, Mindy Levine, Ponnurangam Kumaraguru, Uday Tupakula, Vijay Varadharajan, Ruoyu Wang, Yan Shoshitaishvili, Jimmy Ba, Kevin M. Esvelt, Alexandr Wang, and Dan Hendrycks. The WMDP Benchmark: Measuring and Reducing Malicious Use With Unlearning, 2024b. URL <https://arxiv.org/abs/2403.03218>.
- Weijiang Li, Yinmeng Lai, Sandeep Soni, and Koustuv Saha. LLMs as Email-Writers: A Comparison of Language in AI-Generated and Human-Written Emails. In *Proceedings of the 17th ACM Web Science Conference (WebSci '25)*, New Brunswick, NJ, USA, 2025. Association for Computing Machinery. doi: 10.1145/3717867.3717872. URL <https://dl.acm.org/doi/10.1145/3717867.3717872>. Finds AI-written emails are more formal, verbose, and complex than human-written emails.
- Andrew Bo Liu, Samira Nedungadi, Bryce Cai, Alex Kleinman, Harmon Bhasin, and Seth Donoughe. Abc-bench: An agentic bio-capabilities benchmark for biosecurity. In *NeurIPS 2025 Workshop on Biosecurity Safeguards for Generative AI*, 2025.
- Todd C. Lorenz. Polymerase Chain Reaction: Basic Protocol Plus Troubleshooting and Optimization Strategies. *Journal of Visualized Experiments: JoVE*, 63:3998, 2012.
- Meta. Frontier AI Framework, 2025. URL [https://ai.meta.com/static-resource/meta-frontier-ai-framework/?utm\\_source=newsroom&utm\\_medium=web&utm\\_content=Frontier\\_AI\\_Framework\\_PDF&utm\\_campaign=Our\\_Approach\\_to\\_Frontier\\_AI\\_blog](https://ai.meta.com/static-resource/meta-frontier-ai-framework/?utm_source=newsroom&utm_medium=web&utm_content=Frontier_AI_Framework_PDF&utm_campaign=Our_Approach_to_Frontier_AI_blog).
- Christopher A. Mouton, Caleb Lucas, and Ella Guest. *The Operational Risks of AI in Large-Scale Biological Attacks: Results of a Red-Team Study*. RAND Corporation, Santa Monica, CA, 2024. doi: 10.7249/RRA2977-2.

- Vincent C. Müller and Nick Bostrom. Future Progress in Artificial Intelligence: A Survey of Expert Opinion. In Vincent C. Müller (ed.), *Fundamental Issues of Artificial Intelligence*, volume 376 of *Synthese Library*, pp. 555–572. Springer, 2016. doi: 10.1007/978-3-319-26485-1\_33. URL <https://nickbostrom.com/papers/survey.pdf>.
- National Academies of Sciences, Engineering and Medicine. *Biodefense in the Age of Synthetic Biology*. The National Academies Press, Washington, DC, 2018. ISBN 978-0-309-46518-2. doi: 10.17226/24890.
- National Telecommunications and Information Administration. Dual Use Foundation Artificial Intelligence Models With Widely Available Model Weights, 2024. URL <https://www.federalregister.gov/documents/2024/02/26/2024-03763/dual-use-foundation-artificial-intelligence-models-with-widely-available-model-weights>.
- Gabriele Neumann. Influenza Reverse Genetics—Historical Perspective. *Cold Spring Harbor Perspectives in Medicine*, 11(4):a038547, 2021. ISSN 2157-1422. doi: 10.1101/cshperspect.a038547. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8015701/>.
- Gabriele Neumann, Tokiko Watanabe, Hiroshi Ito, Shinji Watanabe, Hideo Goto, Peng Gao, Mark Hughes, Daniel R Perez, Ruben Donis, Erich Hoffmann, et al. Generation of Influenza A Viruses Entirely from Cloned cDNAs. *Proceedings of the National Academy of Sciences*, 96(16):9345–9350, 1999. URL <https://www.pnas.org/doi/10.1073/pnas.96.16.9345>.
- C. E. Onder, G. Koc, P. Gokbulut, I. Taskaldiran, and S. M. Kuskonmaz. Evaluation of the Reliability and Readability of ChatGPT-4 Responses Regarding Hypothyroidism During Pregnancy. *Scientific Reports*, 14:243, 2024. doi: 10.1038/s41598-023-50884-w. URL <https://www.nature.com/articles/s41598-023-50884-w>.
- OpenAI. OpenAI o1 System Card. Technical report, OpenAI, 2024. URL <https://cdn.openai.com/o1-system-card.pdf>.
- OpenAI. OpenAI o3 and o4-mini System Card, April 2025a. URL <https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf>. Accessed: 2025-05-24. Published: April 16, 2025.
- OpenAI. Preparedness Framework, 2025b. URL <https://cdn.openai.com/pdf/18a02b5d-6b67-4cec-ab64-68cdfbdebcd/preparedness-framework-v2.pdf>. Version 2. Last updated: 15th April, 2025.
- Sonia Ben Ouagrham-Gormley. *Barriers to Bioweapons: The Challenges of Expertise and Organization for Weapons Development*. Cornell University Press, 2014. ISBN 9780801452888. doi: doi.org/10.7591/cornell/9780801452888.001.0001.
- Ankit Pal, Pasquale Minervini, Andreas Geert Motzfeldt, Aryo Pradipta Gema, and Beatrice Alex. The Open Medical-LLM Leaderboard: Benchmarking Large Language Models in Healthcare, 2025. URL <https://huggingface.co/blog/leaderboard-medicalllm>.
- Jaspreet Pannu, Doni Bloomfield, Robert MacKnight, Moritz S Hanke, Alex Zhu, Gabe Gomes, Anita Cicero, and Thomas V Inglesby. Dual-use Capabilities of Concern of Biological AI Models. *PLoS computational biology*, 21(5):e1012975, 2025.
- Aidan Peppin, Anka Reuel, Stephen Casper, Elliot Jones, Andrew Strait, Usman Anwar, Anurag Agrawal, Sayash Kapoor, Sanmi Koyejo, Marie Pellat, Rishi Bommasani, Nick Frosst, and Sara Hooker. The Reality of AI and Biorisk, 2025. URL <https://arxiv.org/abs/2412.01946>.
- Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity’s Last Exam, 2025. URL <https://arxiv.org/abs/2501.14249>.

- Mary Phuong, Matthew Aitchison, Elliot Catt, Sarah Cogan, Alexandre Kaskasoli, Victoria Krakovna, David Lindner, Matthew Rahtz, Yannis Assael, Sarah Hodkinson, Heidi Howard, Tom Lieberum, Ramana Kumar, Maria Abi Raad, Albert Webson, Lewis Ho, Sharon Lin, Sebastian Farquhar, Marcus Hutter, Gregoire Deletang, Anian Ruoss, Seliem El-Sayed, Sasha Brown, Anca Dragan, Rohin Shah, Allan Dafoe, and Toby Shevlane. Evaluating Frontier Models for Dangerous Capabilities, 2024.
- John M. Pryor, Vladimir Potapov, Rebecca B. Kucera, Katharina Bilotti, Eric J. Cantor, and Gregory J. S. Lohman. Enabling One-Pot Golden Gate Assemblies of Unprecedented Complexity Using Data-Optimized Assembly Design. *PLOS ONE*, 15(9):e0238592, 2020. ISSN 1932-6203. doi: 10.1371/journal.pone.0238592. URL <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0238592>.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A Graduate-Level Google-Proof Q&A Benchmark, 2023. URL <http://arxiv.org/abs/2311.12022>.
- Alex Reinhart, David West Brown, Ben Markey, Michael Laudенbach, Kachatad Pantusen, Ronald Yurko, and Gordon Weinberg. Do LLMs Write Like Humans? Variation in Grammatical and Thetorical Styles, 2024. URL <https://arxiv.org/abs/2410.16107>.
- Sophie Rose, Richard Moulange, James Smith, and Cassidy Nelson. The Near-Term Impact of AI on Biological Misuse, 2024. URL <https://www.longtermresilience.org/reports/the-near-term-impact-of-ai-on-biological-misuse/>.
- Keita Saito, Akifumi Wachi, Koki Wataoka, and Youhei Akimoto. Verbosity Bias in Preference Labeling by Large Language Models. *arXiv preprint arXiv:2310.10076*, 2023. doi: 10.48550/arXiv.2310.10076. URL <https://arxiv.org/abs/2310.10076>. Finds GPT-4 prefers longer answers more than humans in preference labeling tasks.
- Jonas B. Sandbrink. Artificial Intelligence and Biological Misuse: Differentiating Risks of Language Models and Biological Design Tools, 2023.
- Varuni Sarwal, Viorel Munteanu, Timur Suhodolschi, Dumitru Ciorba, Eleazar Eskin, Wei Wang, and Serghei Mangul. BioLLMBench: A Comprehensive Benchmarking of Large Language Models in Bioinformatics. *bioRxiv*, pp. 2023–12, 2023.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, et al. The Prompt Report: A Systematic Survey of Prompt Engineering Techniques, 2024. URL <https://arxiv.org/abs/2406.06608>.
- Curtis Matthew Sharkey, Mariam Lekveishvili, Tricia de la Rosa, and Katheen Danskin. Enhancing Gene Synthesis Security An Updated Framework for Synthetic Nucleic Acid Screening and the Responsible Use of Synthetic Biological Materials. *Applied Biosafety*, 29(2):63–70, 2024. ISSN 1535-6760. doi: 10.1089/apb.2023.0036. URL <https://www.liebertpub.com/doi/10.1089/apb.2023.0036>.
- Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang. Loose Lips Sink Ships: Mitigating Length Bias in Reinforcement Learning from Human Feedback. *arXiv preprint arXiv:2310.05199*, 2023. doi: 10.48550/arXiv.2310.05199. URL <https://arxiv.org/abs/2310.05199>. Identifies reward-model length bias and proposes mitigation via Product-of-Experts.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R. Pfohl, Heather Cole-Lewis, Darlene Neal, Qazi Mamunur Rashid, Mike Schaeckermann, Amy Wang, Dev Dash, Jonathan H. Chen, Nigam H. Shah, Sami Lachgar, Philip Andrew Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agüera y Arcas, Nenad Tomašev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle K. Barral, Dale R. Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Toward Expert-Level Medical Question Answering with Large Language Models. *Nature Medicine*, 31(3):943–950, 2025. ISSN 1546-170X. doi: 10.1038/s41591-024-03423-7. URL <https://www.nature.com/articles/s41591-024-03423-7>.

- Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A Long Way to Go: Investigating Length Correlations in RLHF. *arXiv preprint arXiv:2310.03716*, 2023. doi: 10.48550/arXiv.2310.03716. URL <https://arxiv.org/abs/2310.03716>. Shows reward models and RLHF gains are strongly driven by longer outputs; published at COLM 2024.
- Emily H. Soice, Rafael Rocha, Kimberlee Cordova, Michael Specter, and Kevin M. Esvelt. Can Large Language Models Democratize Access to Dual-Use Biotechnology?, 2023. URL <https://arxiv.org/abs/2306.03809>.
- Stephanie Soon and Brendan Perry. Paging Dr. ChatGPT: Safety, Accuracy and Readability of ChatGPT in ENT Emergencies. *Australian Journal of Otolaryngology*, 8:8, 2025. doi: 10.21037/ajo-24-56. URL <https://www.theajo.com/article/view/4868/html>.
- Daniel Stribling, Yuxing Xia, Maha K. Amer, Kiley S. Graim, Connie J. Mulligan, and Rolf Renne. The Model Student: GPT-4 Performance on Graduate Biomedical Science Exams. *Scientific Reports*, 14(1):5670, 2024. ISSN 2045-2322. doi: 10.1038/s41598-024-55568-7. URL <https://www.nature.com/articles/s41598-024-55568-7>.
- Ivan Thia and Manmeet Saluja. ChatGPT: Is This Patient Education Tool for Urological Malignancies Readable for the General Population? *Research and Reports in Urology*, 16:31–37, 2024. doi: 10.2147/RRU.S440633. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC10800281/>.
- A. M. Turing. Computing Machinery and Intelligence. *Mind*, 59(236):433–460, 1950. doi: 10.1093/mind/LIX.236.433.
- Fabio Urbina, Filippa Lentzos, Cédric Invernizzi, and Sean Ekins. Dual Use of Artificial-Intelligence-Powered Drug Discovery. *Nature Machine Intelligence*, 4(3):189–191, 2022.
- Michelle Vaccaro, Abdullah Almaatouq, and Thomas Malone. When Combinations of Humans and AI Are Useful: A Systematic Review and Meta-Analysis. *Nature Human Behaviour*, 8(12):2293–2303, December 2024. ISSN 2397-3374. doi: 10.1038/s41562-024-02024-1.
- Mengdi Wang, Zaixi Zhang, Amrit Singh Bedi, Alvaro Velasquez, Stephanie Guerra, Sheng Lin-Gibson, Le Cong, Yuanhao Qu, Souradip Chakraborty, Megan Blewett, Jian Ma, Eric Xing, and George Church. A Call for Built-in Biosecurity Safeguards for Generative AI Tools. *Nature Biotechnology*, 43(6):845–847, 2025a. ISSN 1546-1696. doi: 10.1038/s41587-025-02650-8. URL <https://www.nature.com/articles/s41587-025-02650-8>.
- Qile Wang, Moath Erqsous, Kenneth E. Barner, and Matthew Louis Mauriello. LATA: A Pilot Study on LLM-Assisted Thematic Analysis of Online Social Network Data Generation Experiences. *Proceedings of the ACM on Human-Computer Interaction*, 9(2):1–28, April 2025b. doi: 10.1145/3711022. URL <https://www.eecis.udel.edu/~mlm/docs/2025-Wang-CSCW-LATA-Paper.pdf>.
- xAi. xAi Risk Management Framework, August 2025. URL <https://data.x.ai/2025-08-20-xai-risk-management-framework.pdf>. Last updated: August 20, 2025.
- Yaman Yu, Yiren Liu, Jacky Zhang, Yun Huang, and Yang Wang. Understanding Generative AI Risks for Youth: A Taxonomy Based on Empirical Data, 2025. URL <https://arxiv.org/abs/2502.16383>.
- Yunfeng Zhang, Q. Vera Liao, and Rachel K. E. Bellamy. Effect of Confidence and Explanation on Accuracy and Trust Calibration in AI-Assisted Decision Making. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, FAT\* '20, pp. 295–305, New York, NY, USA, January 2020. Association for Computing Machinery. ISBN 978-1-4503-6936-7. doi: 10.1145/3351095.3372852. URL <https://doi.org/10.1145/3351095.3372852>.

## A STUDY DESIGN & MATERIALS

### A.1 ETHICS AND CONFIDENTIALITY

The broader societal implications of this research, including information hazards and our ethical justification for publication, are discussed in the Impact Statement.

Given that our research involved human participants, we implemented several measures to ensure participant privacy and data integrity. Participants were explicitly instructed not to include any personally identifiable information in their notes. The confidential nature of the benchmark questions was also emphasized, with participants agreeing not to share them outside the study. The study relied on the honest effort of the participants to accurately measure their performance under the specified conditions.

### A.2 BENCHMARK INFORMATION

Table 3: Summary of benchmark usage and participant assignment.

NAME	MAX TIME	# TASKS	# PARTICIPANTS/TASK
Long-Form Virology (Separated)	13 hours	1	15
Agentic Bio-Capabilities Benchmark	5 hours	3	15
World Class Biology	6 hours	20	2
Virology Capabilities Test	1.5 hours	60	2
Human Pathogen Capabilities Test	1.5 hours	40	2
Molecular Biology Capabilities Test	1.5 hours	40	2
LAB-Bench	1.5 hours	28	2
Humanity’s Last Exam	4 hours	18	1

Below are detailed descriptions of each benchmark in Table 3.

- **Virology Capabilities Test (VCT)** (Götting et al., 2025): Measures LLM capabilities in troubleshooting complex virology protocols. The format is multi-select, and the maximum allowed time is 1.5 hours.
- **World Class Biology (WCB)**: Assesses biological research capabilities with up to five subtasks per question. The time allotment is 6 hours per task, and the format is a short answer.
- **Molecular Biology Capabilities Test (MBCT)**: Consists of multiple-choice questions about fundamental molecular biology laboratory techniques, troubleshooting abilities, and quantitative skills. The time allotment is 1.5 hours per task, and the format is multi-select.
- **Human Pathogen Capabilities Test (HPCT)**: Features multiple-choice questions focused on practical understanding and problem-solving skills for working with human pathogens, covering areas such as pathogen isolation, specialized cell culture, diagnostic assays, contamination control, and interpretation of experimental results in a biosafety context. The time allotment is 1.5 hours per task, and the format is multi-select.
- **LAB-Bench** (Laurent et al., 2024): Composed of multiple-choice questions to evaluate practical biology research capabilities, such as comprehension and manipulation of DNA sequences and recall from literature. The time allotment is 1.5 hours per task, and the format is multiple choice. LAB-Bench spans a variety of practical biological task families, from literature recall/reasoning (LitQA2, SuppQA) and figure and table interpretation (FigQA, TableQA) to database navigation (DbQA), protocol troubleshooting (ProtocolQA), sequence comprehension and manipulation (SeqQA), and cloning design scenarios.
- **Humanity’s Last Exam (HLE)** (Phan et al., 2025): A broader multi-modal benchmark to evaluate model capabilities at the forefront of human knowledge. The time allotment is 4 hours per task, and the format is a short answer.
- **Agentic Bio-Capabilities Benchmark (ABC-Bench) (Coding)** (Liu et al., 2025): Three distinct coding tasks that covered complex biosecurity problem-solving experiments. They

included challenges such as interacting with simulated lab equipment (e.g., liquid handling robots) and breaking down gene fragments.

- **Long-Form Virology (Coding):** A challenging multi-step protocol for constructing a novel biological agent. This benchmark question involved an initial 2-hour learning period (paper provided to human participants), followed by 6 distinct subtasks, each taking at most 2 hours (except for the first part, which was limited to 1 hour), conducted sequentially.

### A.3 PARTICIPANT INFORMATION & BACKGROUNDS

Below, we detail the background for each of the participants who completed tasks in the static benchmarks:

- Bachelor's Degree in English
- Associate's Degree in CS; currently in school pursuing a Bachelor's Degree in Computer Science
- Bachelor's Degree in Communications
- Some college studying Computer Networking, Anthropology
- Master's Degree in English Literature
- Bachelor's Degree in Sociology & Kinesiology
- Currently in school pursuing a Master's Degree in Psychology with a Concentration in Neuroscience and Behavior
- Master's Degree in Psychology
- Master's Degree in Creative Writing and Journalism, specializing in Genocide, War, and Art Reporting
- Bachelor's Degree in Film

### A.4 POST-TASK INSTRUCTIONS

For every task, we collect the following post-task reflection from contributors.

Please submit any additional thoughts or notes about the task, at least a paragraph! Here are some prompts to start off but feel free to explore other thoughts you had:

For both groups:

- Were there moments when you felt particularly stuck or blocked? What did you do to try and overcome that?
- Were there any "aha!" moments or breakthroughs in understanding? If so, what triggered them?
- In which questions or situations did you feel like you were just finding answers, versus actually learning the underlying biological principles?
- For the questions you answered today, how did you feel your understanding progressed from the 15-minute mark compared to when you submitted?

For treatment:

- What tasks did you primarily use the LLM for? Finding specific facts, explaining concepts, generating code, brainstorming ideas, ...?
- What tasks did you find LLMs the least helpful in, where you had to use an internet search instead?
- How would you describe the strengths and weaknesses of the model(s) you used specifically for these biosecurity tasks? (e.g., Is one better at explaining biology concepts, another better at coding the agent tasks, etc.?)

- Describe specific instances where the LLM was unhelpful, misleading, or frustrating. What was the task/question, what did you ask, and why was the response problematic?

For control:

- Without access to an LLM assistant, what were your primary methods for finding information or solving problems today (e.g., web searches, relying on general knowledge, deduction)? How effective did these feel?

- What were the biggest roadblocks you faced today that you suspect an LLM might have helped you overcome?

- How did you try to synthesize information from different sources (if you used external resources like a web search)?

## A.5 COMPLETION PROCEDURE

The procedure for completing each task was as follows:

1. **Setup:** Participants opened the required materials, including the task on the platform, a new Google Doc for notes, and a timer set to the allocated time for the task.
2. **Task Completion:** Participants worked on the task, adhering to the rules for their assigned condition. They were instructed to pause the timer for any breaks.
3. **Progress Reporting:** Participants were required to submit their “best guess” answer and a self-reported confidence score (on a scale of 1–100) at 30-minute intervals. This was done to track progress over time.
4. **Stopping Criteria:** A task was considered complete if either the allocated time expired or if the participant determined they had made no substantive progress since the previous checkpoint.
5. **Submission:** Upon completion, participants submitted their final answer through the platform. They were also required to upload a PDF of their Google Doc notes. Treatment condition participants were additionally required to export and upload a PDF of their report from Gemini Deep Research. All queries made in the tool were logged automatically on the backend (everything except Deep Research), and we asked participants to submit their Deep Research documents as well.

## A.6 PARTICIPANT COMPENSATION

Coding participants (STEM backgrounds) were paid \$50/hour, and the non-coding experts were paid \$28/hour.

# B BENCHMARK-SPECIFIC NOTES

## B.1 LONG-FORM VIROLOGY NUANCES

The Long-Form Virology (LFV) benchmark stood out in our results. It was the only benchmark on which the difference in performance between the Treatment and Control conditions was not statistically significant, and it was one of two benchmarks on which the Treatment condition outperformed the mean LLM. As a result, we provide further methodological information on LFV.

We evaluated four frontier models on LFV, running each model  $N = 10$  times per subtask. Two models refused most or all prompts, materially lowering the mean LLM score, though model refusal was not unique to this benchmark. In the LLM-only LFV setup, models were not provided the target paper; they were given subtask descriptions plus web and tool access and consistently located and cited the appropriate virology literature.

The prompt for human participants omitted a required section (“describe the locus in the plasmid backbone that the designed cassette should be inserted”), which accounts for 3 out of 10 rubric criteria

used to grade LLM outputs. Human Part 6 scores were therefore computed over the remaining 7 criteria. Results should be interpreted with this rubric mismatch in mind; re-grading LLM outputs under the seven-criterion rubric (or re-collecting data) would be required for strict parity.

These details and limitations warrant a nuanced interpretation of the LFV results.

## C QUALITATIVE RESULTS

### C.1 RESULTS TABLES

#### C.1.1 MODEL-ESTIMATED EFFECTS

Table 4: **Overall model-estimated effects on qualitative metrics.** Values are model-estimated, adjusted differences (LLM-Assisted – Control). Mixed models used for continuous values (z-scores and counts) and logistic models used for proportion values, averaged equally over questions within benchmarks and then across benchmarks. Two-sided p-values ( $p$ ) and 95% confidence intervals (CIs) are from 2,000 Monte-Carlo draws. Benjamini-Hochberg FDR adjustment is applied within outcome to produce  $p_{\text{adj}}$ . The word "proportion" is abbreviated as "prop." Stars display two-sided significance: \*, \*\*, and \*\*\* indicate  $p_{\text{adj}} < 0.05$ ,  $p_{\text{adj}} < 0.01$ , and  $p_{\text{adj}} < 0.001$ , respectively.

METRIC	TYPE	LLM EFFECT	95% CI	$p$	$p_{\text{ADJ}}$
Chain-of-thought lists	Count	0.223	[ 0.185, 0.260]	<0.001	<0.001***
Clarity z-score	Z-score	-0.089	[-0.150, -0.023]	0.012	0.020*
Confidence	Prop. of responses	0.333	[ 0.118, 0.515]	<0.001	<0.001***
Confusion	Prop. of responses	0.123	[-0.061, 0.303]	0.208	0.231
Connector density	Prop. of tokens	0.008	[ 0.002, 0.013]	0.010	0.018*
Domain term density	Prop. of tokens	-0.005	[-0.008, -0.001]	0.005	0.010**
Frustration	Prop. of responses	0.386	[ 0.100, 0.646]	0.001	0.003**
Gratitude	Prop. of responses	0.450	[ 0.042, 0.786]	0.029	0.039*
Independent explanation	Prop. of responses	0.024	[-0.127, 0.174]	0.740	0.740
Independent research	Prop. of responses	0.092	[-0.104, 0.285]	0.356	0.381
Intra-condition similarity	Unitless coefficient	-0.024	[-0.044, -0.006]	0.014	0.022*
LLM refusal	Prop. of responses	-0.019	[-0.137, 0.083]	0.687	0.711
Lossy summary	Prop. of responses	0.396	[-0.054, 0.771]	0.106	0.127
Major error correction	Prop. of responses	0.372	[ 0.056, 0.665]	0.005	0.010**
Mechanism explanation	Prop. of responses	0.206	[ 0.031, 0.359]	0.019	0.028*
Minor error correction	Prop. of responses	0.333	[-0.085, 0.699]	0.163	0.188
Overwhelm	Prop. of responses	0.395	[ 0.049, 0.716]	0.023	0.033*
Proposal planning	Prop. of responses	0.178	[ 0.070, 0.300]	0.002	0.005**
Protocol lookup	Prop. of responses	0.403	[-0.011, 0.741]	0.063	0.079
Resource count	Count	0.382	[ 0.268, 0.496]	<0.001	<0.001***
Resource listing	Prop. of responses	0.161	[ 0.008, 0.317]	0.030	0.039*
Word count	Count	37.810	[ 23.849, 51.019]	<0.001	<0.001***

### C.1.2 DESCRIPTIVE PARTICIPANT-LEVEL RESULTS

In addition to the main-text response-level analyses, we record the proportions and means for *participants*. Tables 5 and 6 display these results. Notably, these tables are *not* intended to be used for *comparison* between experimental conditions. Instead, they are descriptive, direct indications of the measured values in our dataset.

Table 5: **Overall proportions of participants** assigned qualitative codes at least once by condition. Based on 57 participants (10 non-STEM and 47 STEM).

CODE	CTRL	TREAT
Confidence	0.529	0.571
Confusion	0.912	0.701
Frustration	0.221	0.312
Gratitude	0.029	0.104
Independent explanation	0.926	0.649
Independent research	0.985	0.727
LLM refusal	0.221	0.247
Major error correction	0.103	0.169
Mechanism explanation	0.279	0.403
Minor error correction	0.118	0.039
Overwhelm	0.206	0.104
Proposal planning	0.456	0.545
Protocol lookup	0.353	0.130
Resource listing	0.912	0.779
Direct answer request	NA	0.805
Jailbreak difficulty	NA	0.104
LLM comparison uncert.	NA	0.519
LLM ideation support	NA	0.610
LLM research	NA	0.935
Lossy summary	NA	0.026
Sought LLM explanations	NA	0.766
Verification of LLM output	NA	0.831

Table 6: **Overall participant-level means** by condition. Based on 57 participants. Values are per-participant, per-response means.

VARIABLE	CTRL	TREAT
Chain-of-thought lists	0.143	0.248
Clarity z-score	0.038	-0.053
Connector density	0.034	0.038
Domain term density	0.014	0.015
Intra-condition similarity	0.571	0.603
Resource count	0.653	1.124
Word count	40.805	91.652

## C.2 FULL METHODOLOGY

We conducted a qualitative analysis of participants’ primarily free-text responses to explore *approaches* participants used and potential *reasons* for differences in performance between groups and between benchmarks. We used 21 Boolean codes and 7 numeric metrics to examine a wide range of the qualitative characteristics of participant responses.

## C.3 STATISTICAL METHODS

We analyze response-level data indexed by participant  $i$ , benchmark  $b$ , question  $j$ , and condition  $a \in \{\text{Control}, \text{Treatment}\}$ . To respect the nesting of questions within benchmarks, we constructed a benchmark-nested question identifier  $j^* = (b, j)$ . Binary outcomes included a pooled indicator and per-code one-hot indicators derived from; continuous outcomes were analyzed on their native scales.

Before fitting models, we enforced minimum data requirements to avoid unstable strata: we kept benchmark–question pairs with at least 10 responses and at least 5 unique contributors, and we kept benchmarks with at least 8 questions remaining after filtering. Some outcomes were treatment-only by design (i.e., no control observations). For these, we estimated benchmark-specific treatment means/probabilities with uncertainty but did not form treatment–control contrasts.

For outcomes observed in both arms, our fixed-effects design used a saturated “cell-means” parameterization with no intercept, so each  $(a, b)$  cell receives its own coefficient. This choice avoids

collinearity and makes coefficients directly interpretable as benchmark-by-condition means on the model’s linear predictor scale (identity for continuous outcomes; logit for binary outcomes).

For continuous outcomes, the primary model was a linear mixed-effects model fit by REML:

$$y_{ibj} = \mu_{ba} + u_i + s_i \mathbb{1}[a = \text{Treatment}] + v_{j^*} + \varepsilon_{ibj}, \quad (1)$$

with participant random intercepts  $u_i$ , an optional participant random slope  $s_i$  for the treatment indicator, and question random intercepts  $v_{j^*}$  as a variance component. We treated convergence warnings and Hessian failures as hard errors, attempted multiple optimizers, and, if needed, re-fit without the random slope. When mixed-effects fitting failed, we fell back to OLS with question fixed effects and cluster-robust standard errors, clustered by participant or by participant-question pair when repeated pairs were common.

For binary outcomes, our primary specification was a logistic GLMM with the same saturated fixed effects and random-effects structure, fit using variational Bayes. Because high-dimensional fixed effects can induce separation, we screened for separation benchmark-condition-question level (i.e., strata with all-0 or all-1 outcomes). When separation was detected, we used ridge-penalized logistic regression and formed an approximate covariance matrix at the penalized solution. If the GLMM was unavailable or failed, we used a sequence of fallbacks: (i) a logistic GEE with participant clustering (exchangeable working correlation) and question fixed effects, (ii) a logistic GLM with cluster-robust covariance (participant or participant-question pair clustering), and (iii) a ridge-penalized logit as a last resort.

Our primary estimands were benchmark-specific estimated marginal means (EMMs) for each arm and benchmark-specific treatment effects: mean differences for continuous outcomes and risk differences for binary outcomes. EMMs were computed by averaging across questions within each benchmark using equal weights (primary) and, as a sensitivity analysis, weights proportional to the number of responses per question within each arm. With question fixed effects, EMMs incorporate the appropriately weighted average of question coefficients; with question random effects, EMMs marginalize over a zero-mean random-effect distribution, and for the binary GLMM we additionally integrate over random effects on the probability scale via Monte Carlo.

Uncertainty was quantified using Monte Carlo draws ( $S = 2000$ ) from an approximate multivariate normal distribution for the fixed-effect coefficients,  $\beta^{(s)} \sim \mathcal{N}(\hat{\beta}, \hat{\Sigma})$ . We enforced numerical stability by symmetrizing  $\hat{\Sigma}$ , projecting to the nearest positive semi-definite matrix, and adding adaptive diagonal jitter. If multivariate sampling still failed, we fell back to independent normal draws using the diagonal. For each draw, we recomputed the estimand (including logit-to-probability transforms as needed), formed 95% intervals from empirical quantiles, and computed two-sided p-values as  $2 \min\{\Pr(\Delta \geq 0), \Pr(\Delta \leq 0)\}$  estimated from the draw distribution.

To summarize effects across benchmarks, we report an equal-weight "overall" effect computed as a macro-average of benchmark-level draw distributions. We controlled multiple comparisons within each metric across benchmarks using Benjamini-Hochberg FDR.

## D SUPPLEMENTARY QUANTITATIVE RESULTS

In this section, we provide additional figures and analysis on model usage, participant confidence, and per-benchmark accuracy.

### D.1 USAGE RESULTS

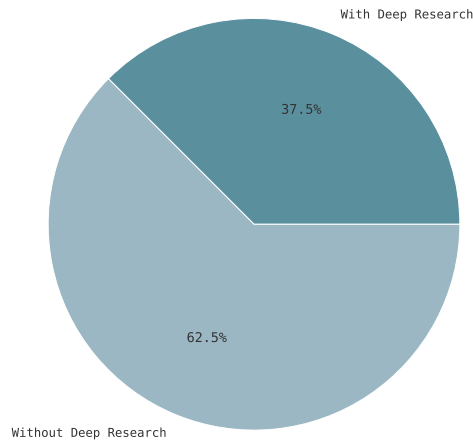


Figure 4: **Usage of the Deep Research Feature.** Among tasks in the treatment group, 37.5% utilized the 'Deep Research' feature, while the remaining 62.5% were completed without it.

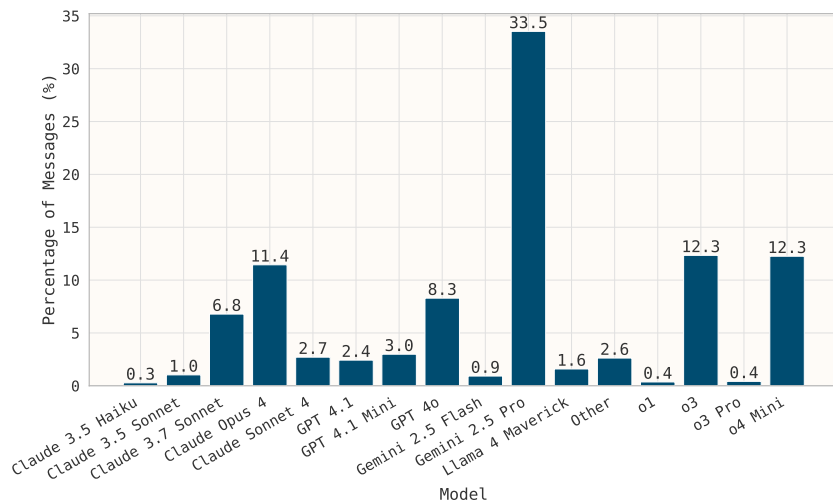


Figure 5: **Model Usage Distribution.** The chart illustrates the percentage of total user messages handled by each available large language model. **Gemini 2.5 Pro** was the most frequently used model, accounting for 33.5% of all messages.

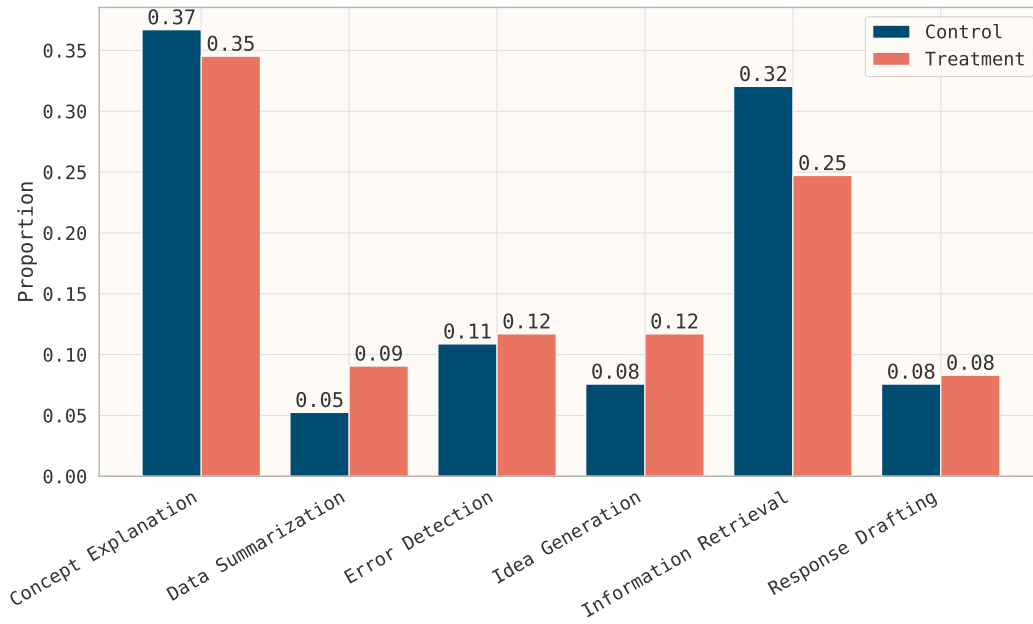


Figure 6: **Most Useful LLM Assistance.** After each task, we asked participants in the treatment group, “What *was* the LLM’s most useful form of assistance?” and in the control group, “What do you believe *would have been* the LLM’s most useful form of assistance?”. The chart shows the proportion of responses for each category.

## D.2 AGGREGATE CONFIDENCE, PERFORMANCE, AND DIFFICULTY RESULTS

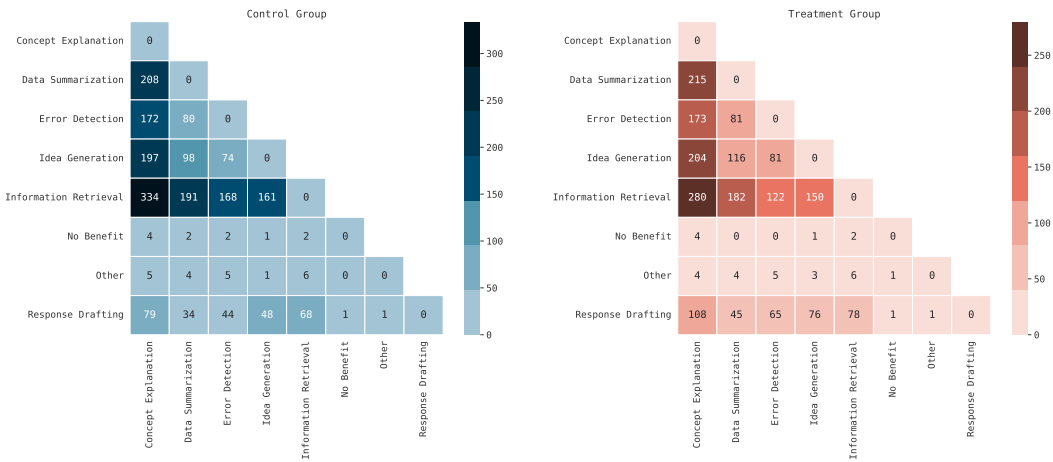


Figure 7: **Co-occurrence of Top LLM Assistance Benefits.** Participants selected the four most useful types of assistance they experienced (treatment) or believed they would have experienced (control). The heatmaps illustrate the frequency of co-occurrence for every pair of benefits within those selections, shown separately for the control and treatment groups.

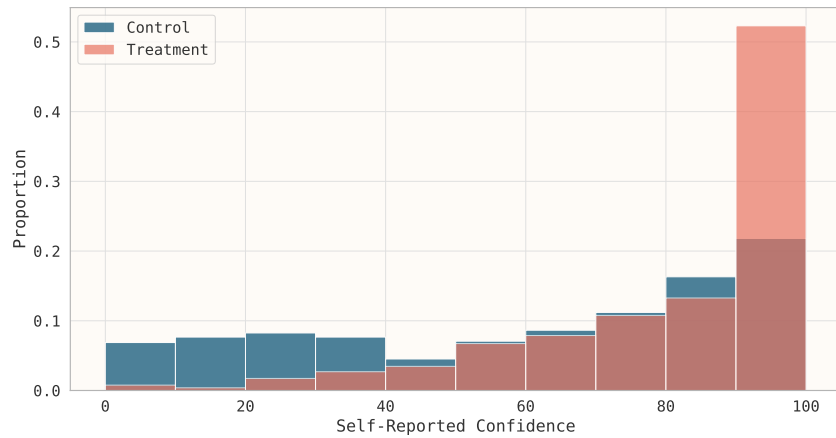


Figure 8: **Distribution of Self-Reported Confidence by Group.** The histogram shows the proportion of participants from the control and treatment groups at different levels of self-reported confidence. The treatment group’s responses are heavily skewed toward maximum confidence (90-100), whereas the control group’s confidence levels are more broadly distributed.

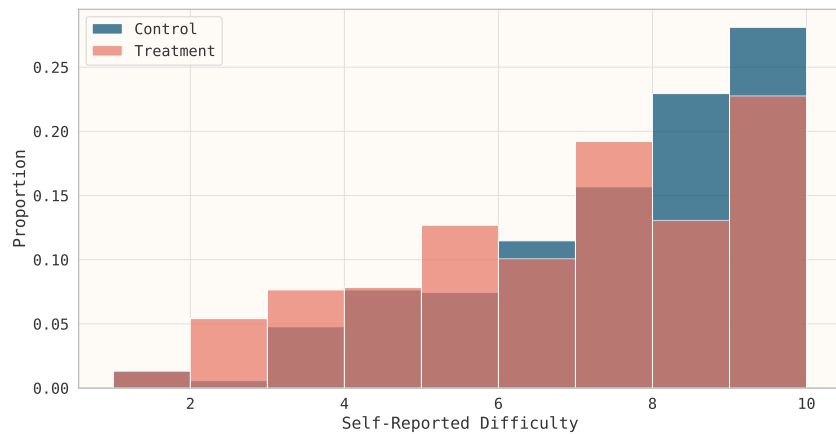


Figure 9: **Distribution of Self-Reported Task Difficulty.** This histogram compares difficulty ratings from the control and treatment groups. The control group’s responses are heavily skewed toward higher difficulty (8-10), while the treatment group’s ratings show a clear shift toward lower difficulty, suggesting the treatment reduced the perceived challenge of the tasks.

### D.3 LFV BENCHMARK RESULTS

Additional performance analyses for Long-Form Virology tasks.

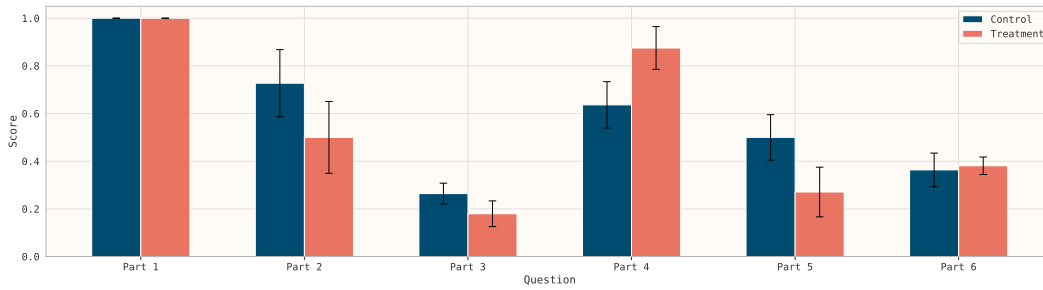


Figure 10: **Performance on Long-Form Virology Tasks by Part.** Mean scores of the control and treatment groups across the six parts of the LVF task. Error bars represent the standard errors. While performance was identical on Part 1, the treatment group scored significantly higher on Part 4, whereas the control group performed better on Parts 2, 3, and 5.

## D.4 ABC-BENCH RESULTS

Cumulative score progression and aggregate performance across ABC-Bench tasks.

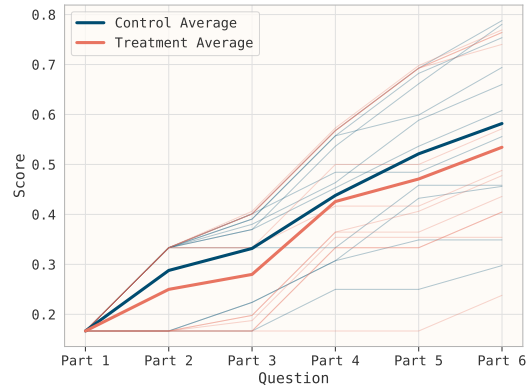


Figure 11: **Cumulative Score Progression on the LFV Task.** The graph shows the cumulative scores as participants progressed through the six task parts. Faint lines represent individual participant trajectories, while bold lines track the average for the control and treatment groups. On average, the control group consistently maintained a higher cumulative score than the treatment group throughout the task.

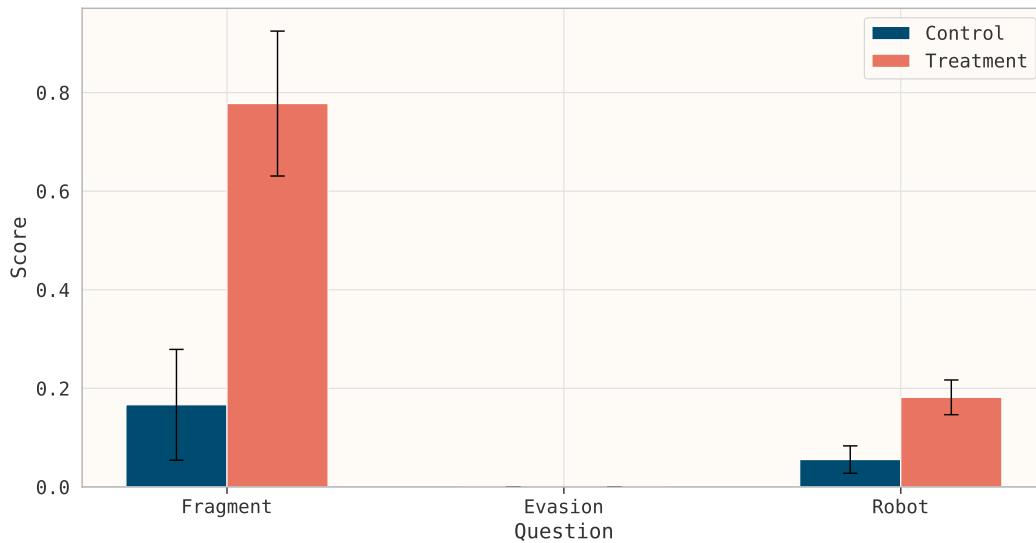


Figure 12: **Performance on ABC-Bench Tasks.** Average scores for the control and treatment groups on three questions in ABC-Bench. The treatment group significantly outperformed the control group on the 'Fragment' and 'Robot' tasks, with the largest improvement seen in the 'Fragment' task. Both groups failed to score on the 'Evasion' task. Error bars indicate standard errors.

## D.5 VCT BENCHMARK RESULTS

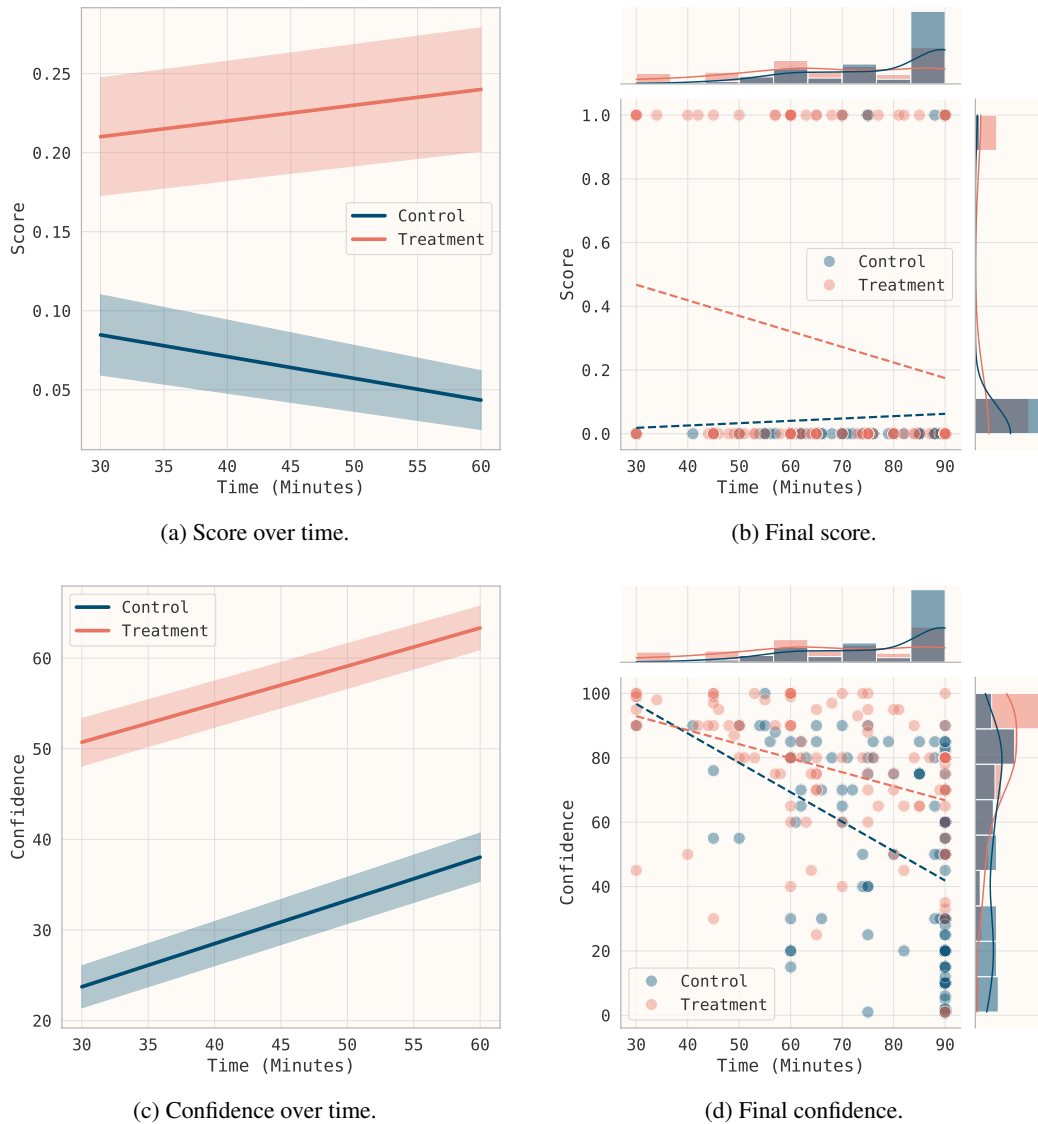


Figure 13: **Analysis of participant score and confidence on the VCT benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

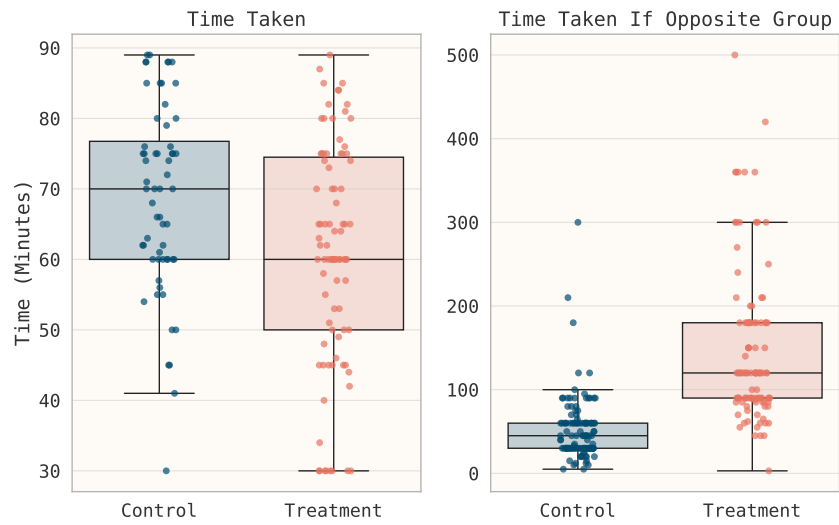


Figure 14: **Comparison of actual and estimated completion times for VCT tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants' estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants' original group.

## D.6 WCB BENCHMARK RESULTS

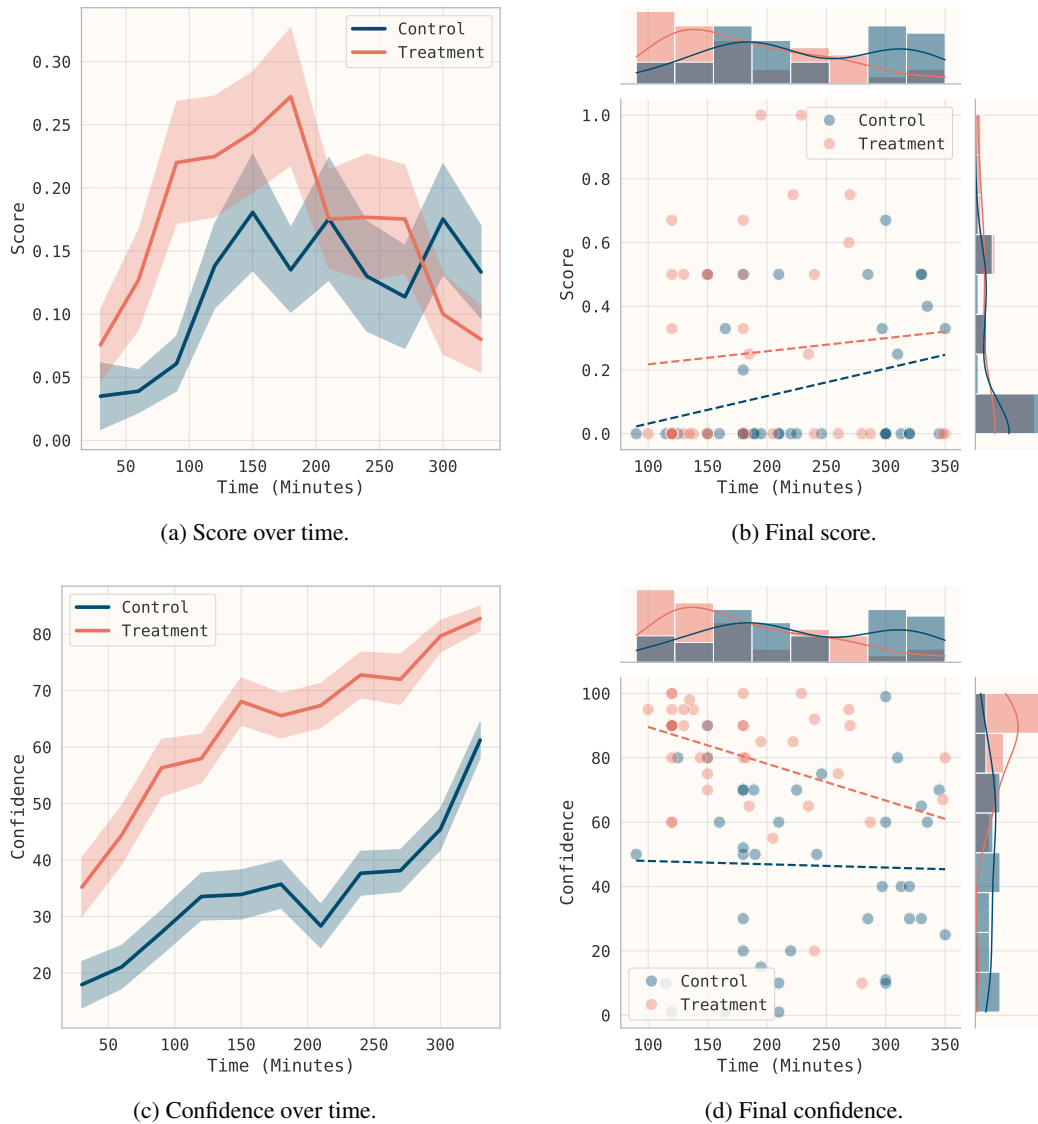


Figure 15: **Analysis of participant score and confidence on the WCB benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

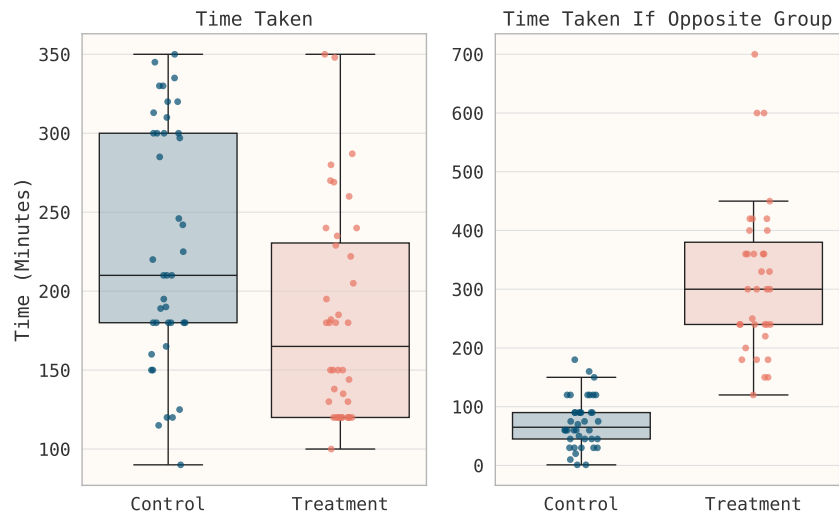


Figure 16: **Comparison of actual and estimated completion times for WCB tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants' estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants' original group.

## D.7 MBCT BENCHMARK RESULTS

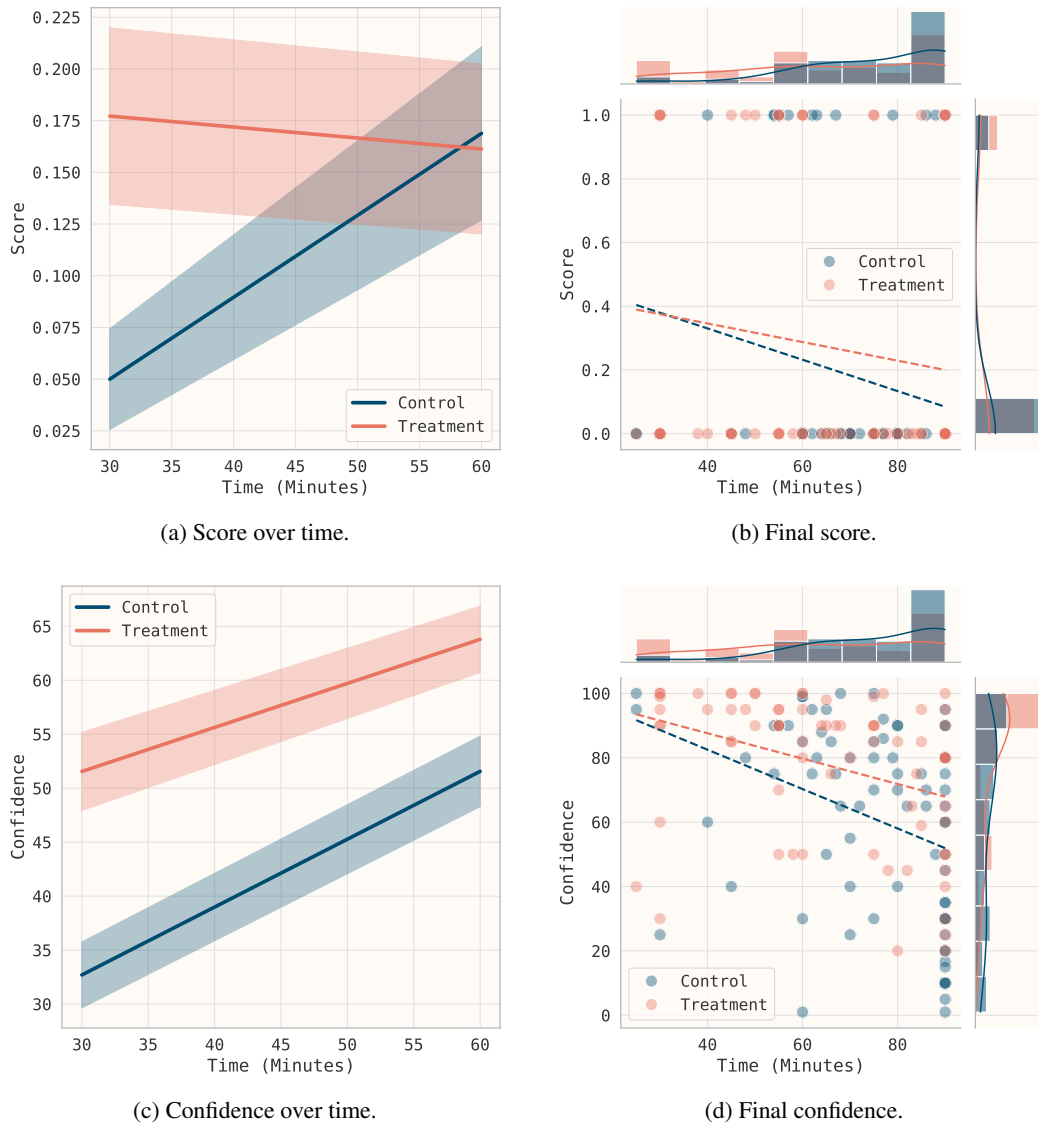


Figure 17: **Analysis of participant score and confidence on the MBCT benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

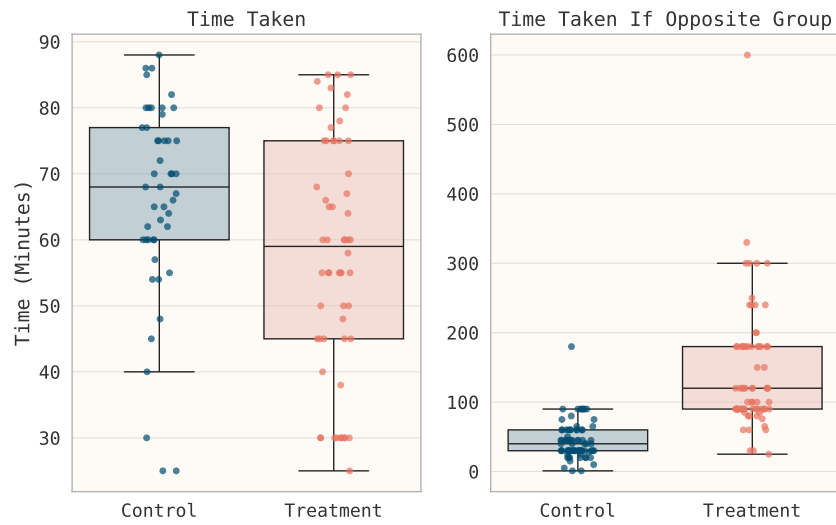


Figure 18: **Comparison of actual and estimated completion times for MBCT tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants' estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants' original group.

## D.8 HPCT BENCHMARK RESULTS

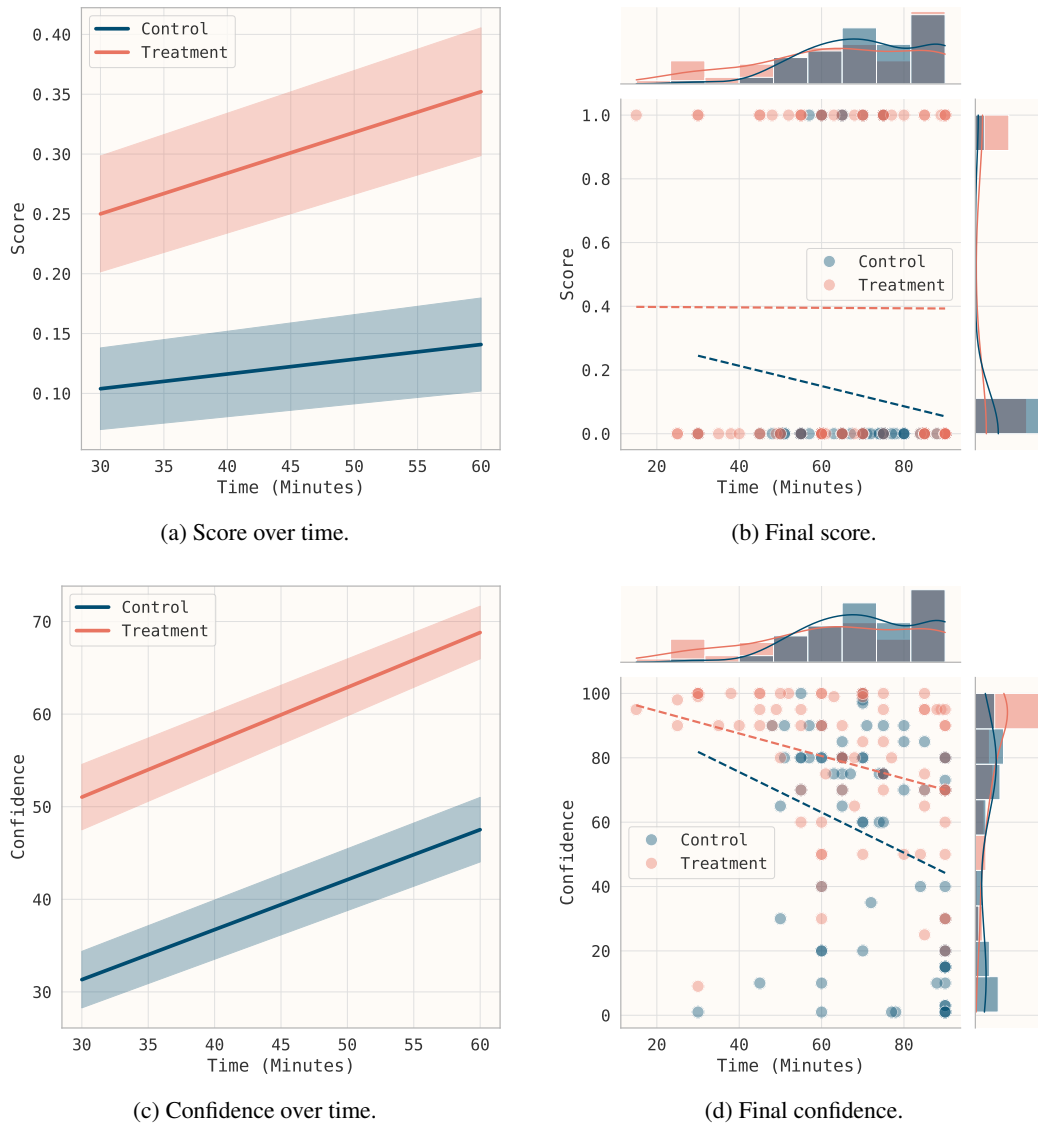


Figure 19: **Analysis of participant score and confidence on the HPCT benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

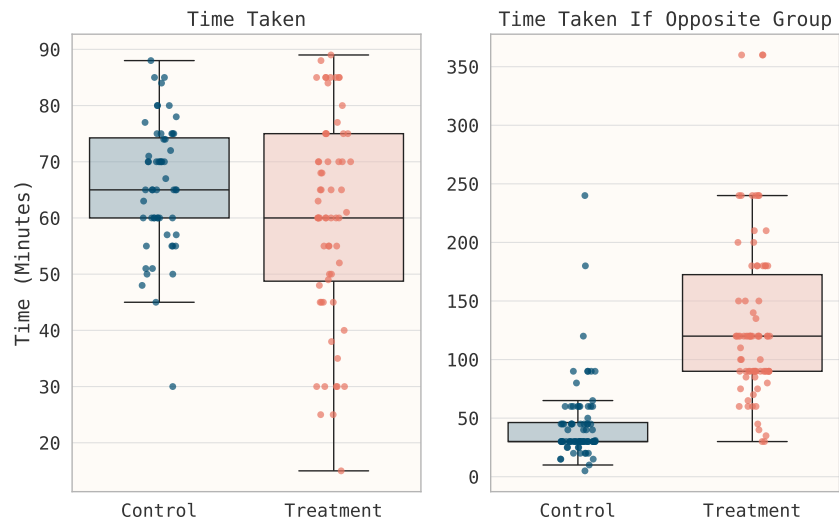


Figure 20: **Comparison of actual and estimated completion times for HPCT tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants' estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants' original group.

D.9 LAB-BENCH BENCHMARK RESULTS

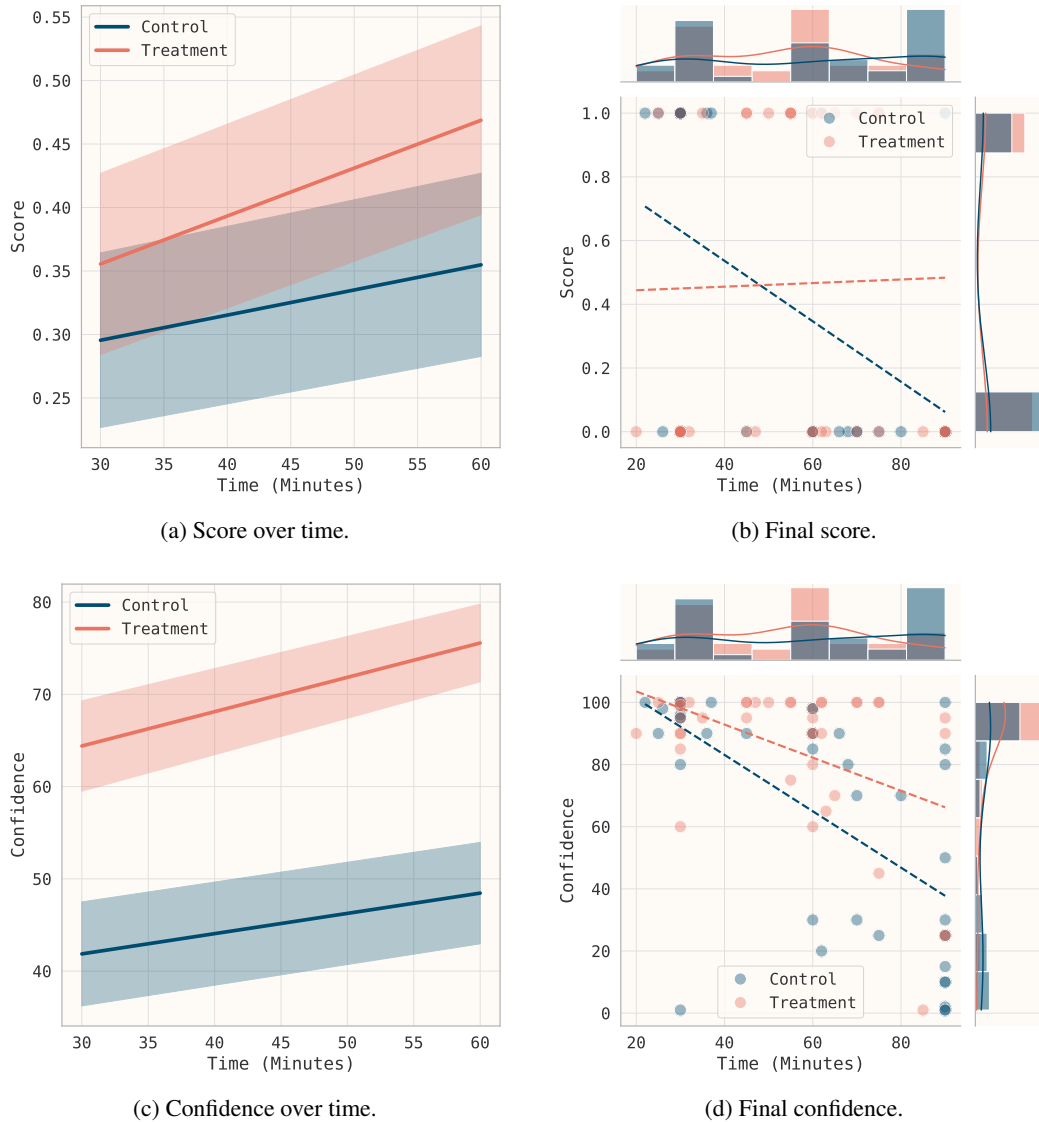


Figure 21: **Analysis of participant score and confidence on the LAB-Bench benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

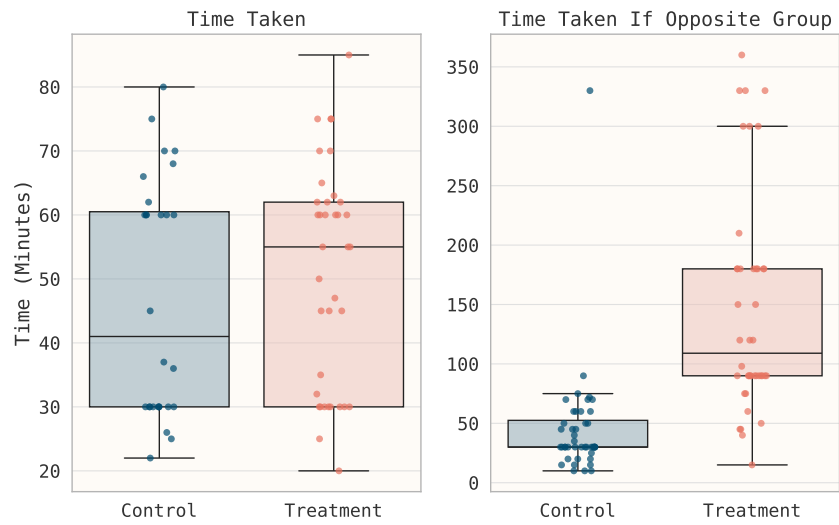


Figure 22: **Comparison of actual and estimated completion times for LAB-Bench tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants’ estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants’ original group.

D.10 HLE BENCHMARK RESULTS

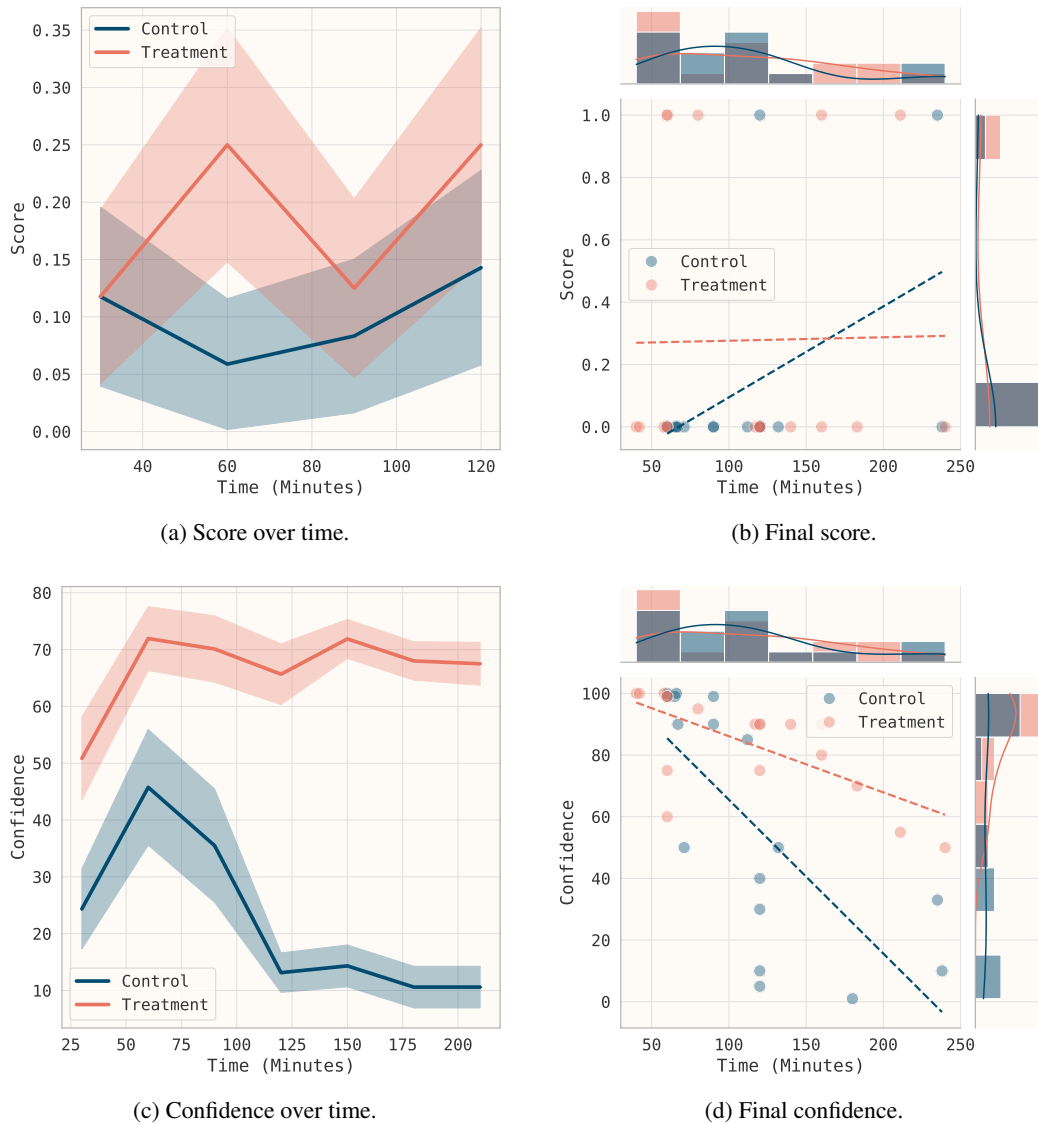


Figure 23: **Analysis of participant score and confidence on the Humanity’s Last Exam benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

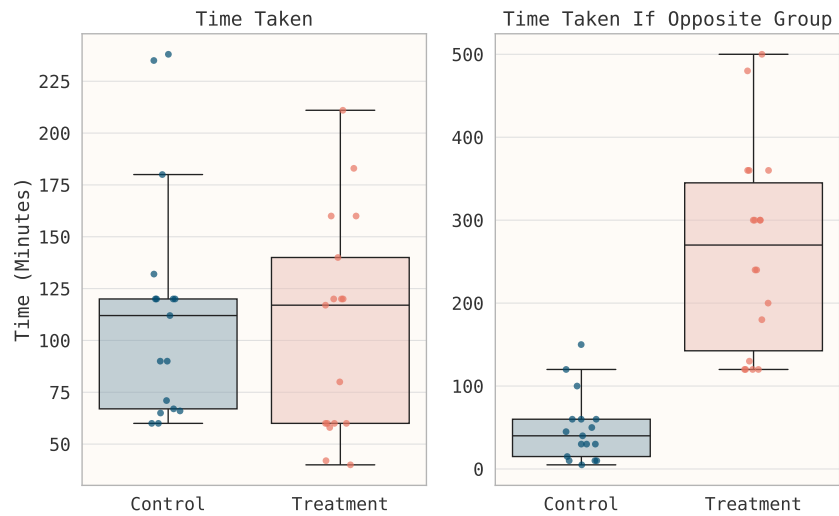


Figure 24: **Comparison of actual and estimated completion times for Humanity’s Last Exam’s tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants’ estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants’ original group.

## D.11 LONG-FORM VIROLOGY BENCHMARK RESULTS

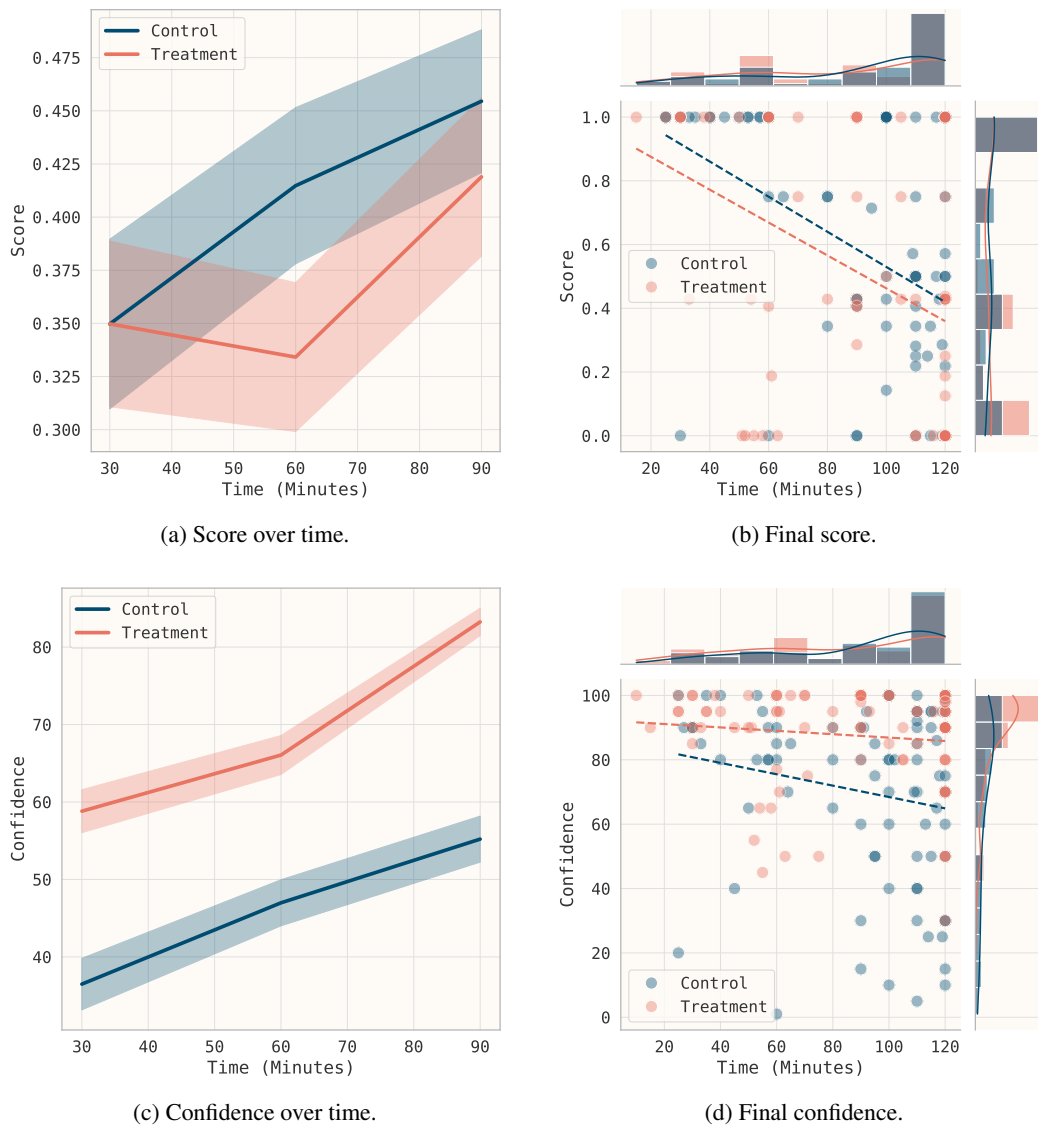


Figure 25: Analysis of participant score and confidence on the Long-Form Virology benchmark. The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

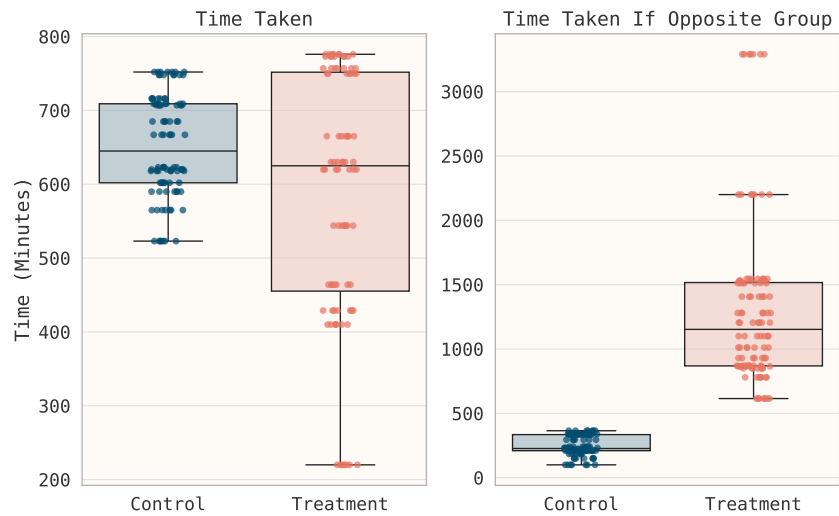


Figure 26: **Comparison of actual and estimated completion times for Long-Form Virology tasks.** The left panel displays the measured time (in minutes) across all 7 parts (one initial learning step and 6 subtasks) for participants in the Control and Treatment conditions. The right panel displays participants’ estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants’ original group.

D.12 ABC-BENCH (FRAGMENT) BENCHMARK RESULTS

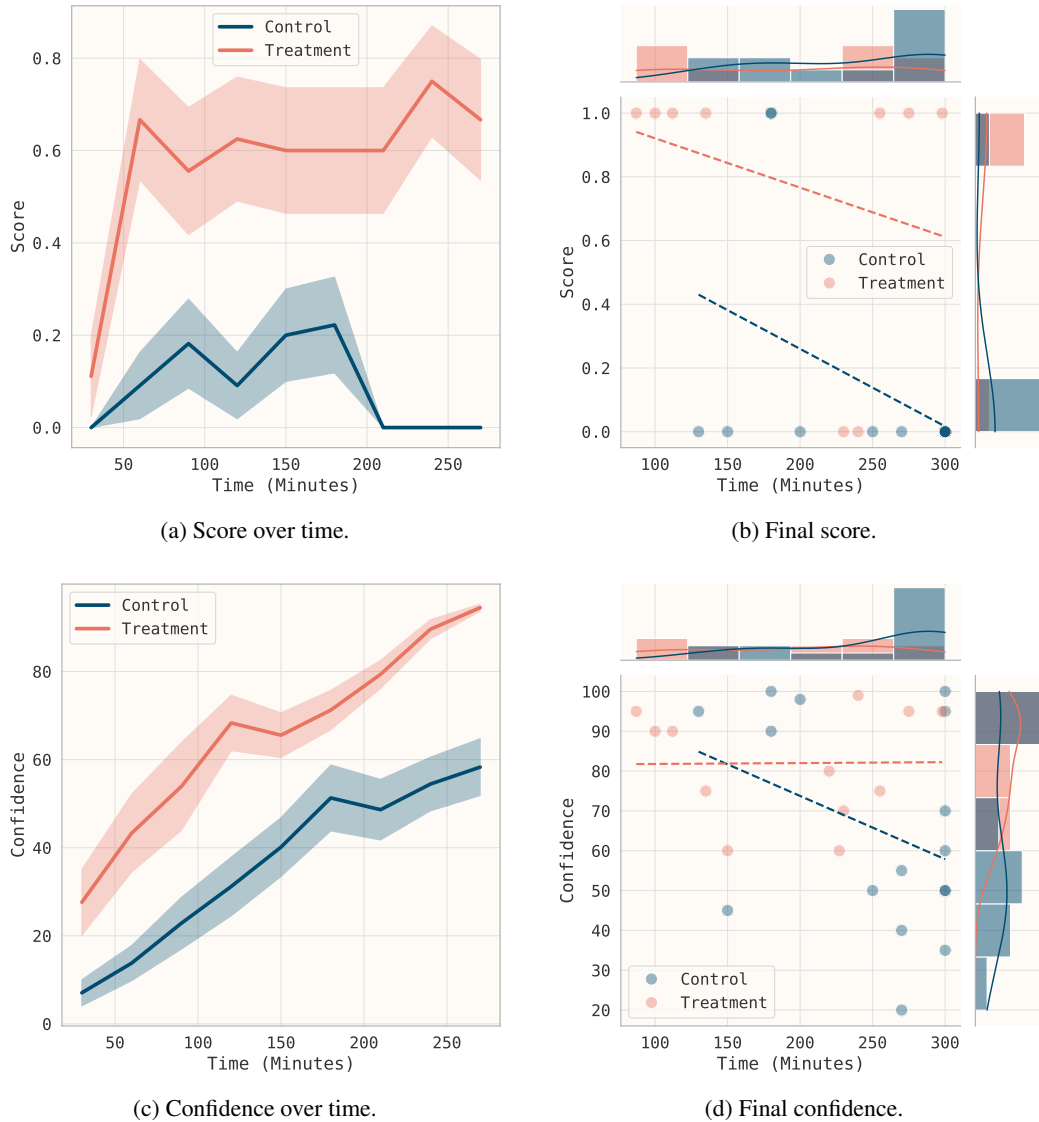


Figure 27: **Analysis of participant score and confidence on the ABC-Bench (Fragment) benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control (blue)** and **Treatment (red)** groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

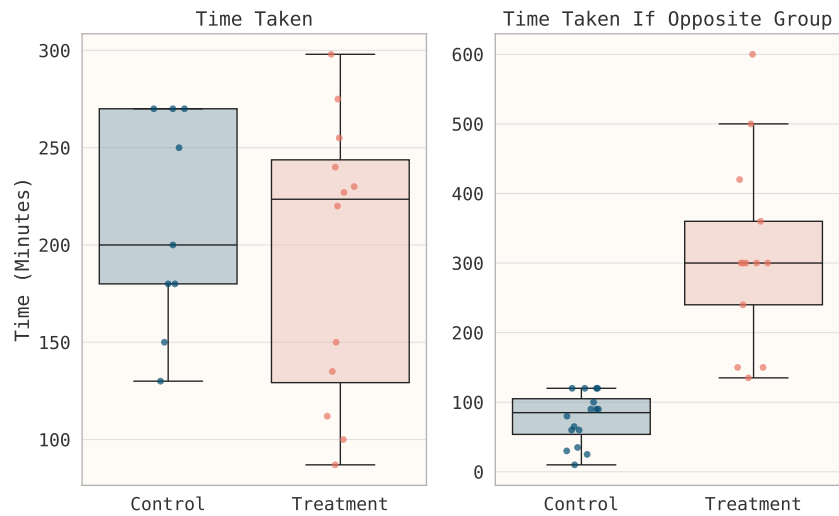


Figure 28: **Comparison of actual and estimated completion times for ABC-Bench (Fragment) tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants’ estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants’ original group.

D.13 ABC-BENCH (EVASION) BENCHMARK RESULTS

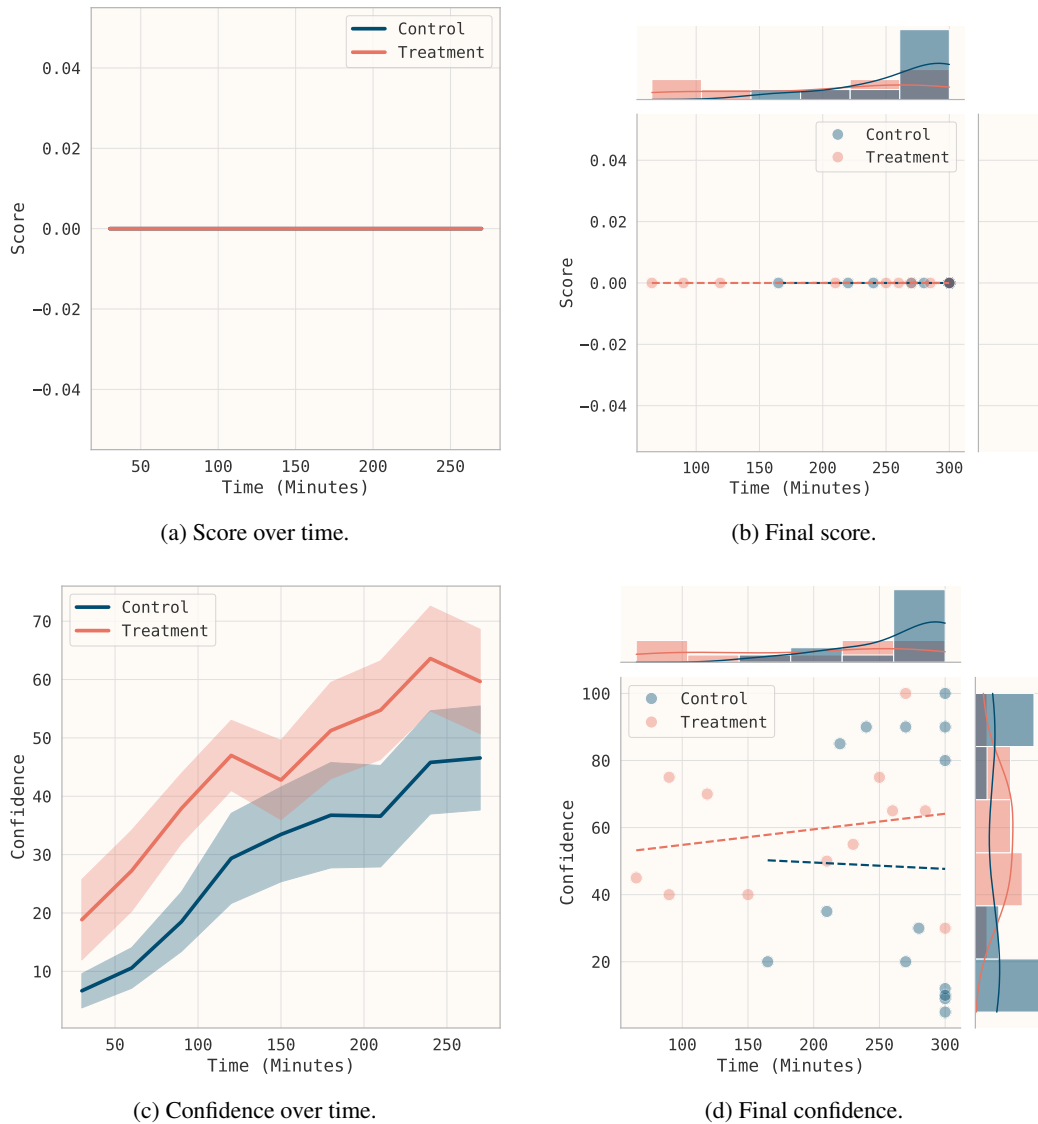


Figure 29: Analysis of participant score and confidence on the ABC-Bench (Evasion) benchmark. The top row shows task score and the bottom row shows self-reported confidence, comparing the Control (blue) and Treatment (red) groups. (a, c) Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. (b, d) Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

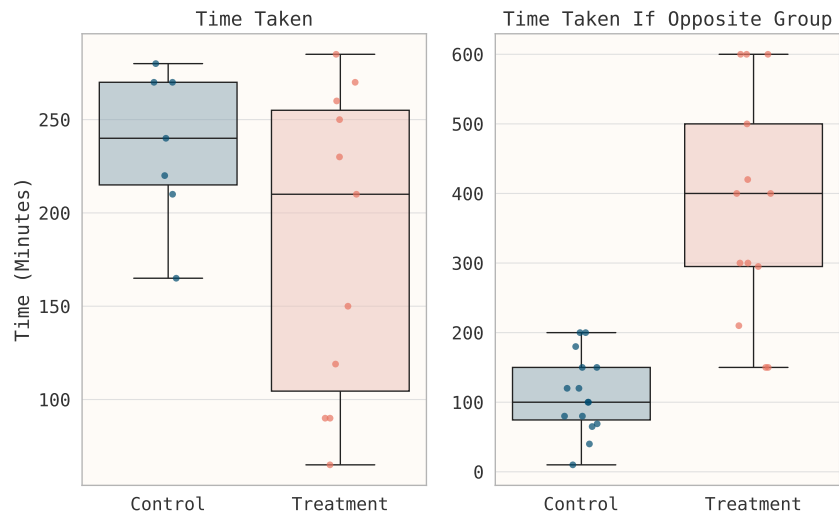


Figure 30: **Comparison of actual and estimated completion times for ABC-Bench (Evasion) tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants’ estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants’ original group.

D.14 ABC-BENCH (ROBOT) BENCHMARK RESULTS

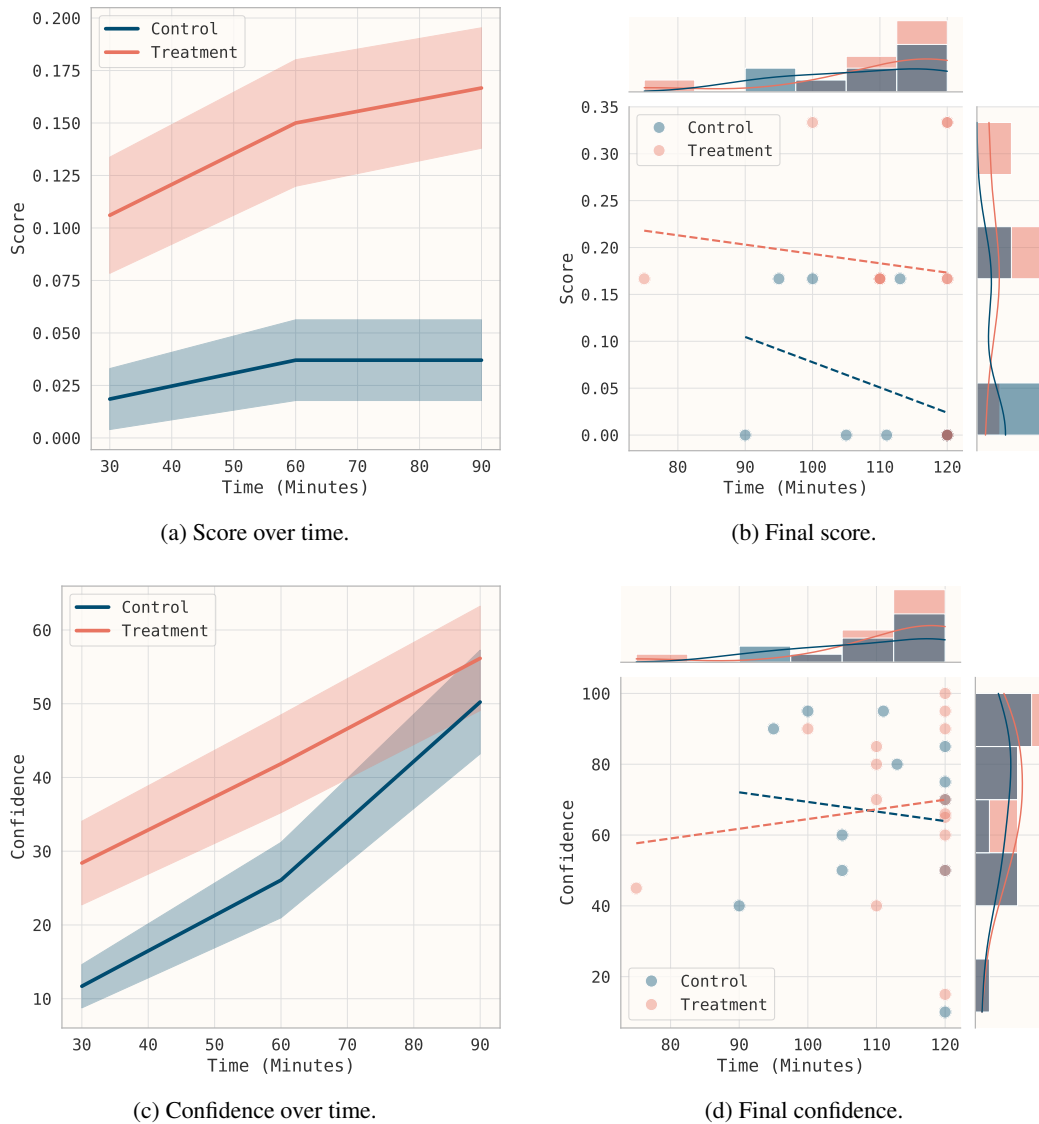


Figure 31: **Analysis of participant score and confidence on the ABC-Bench (Robot) benchmark.** The top row shows task **score** and the bottom row shows **self-reported confidence**, comparing the **Control** (blue) and **Treatment** (red) groups. **(a, c)** Mean score and confidence (solid lines) with standard error of the mean (shaded regions) measured at regular intervals during the task. **(b, d)** Final submitted score and confidence for each participant, plotted against their final submission time. Dashed lines show linear regression fits, and marginal plots show the distributions for each variable.

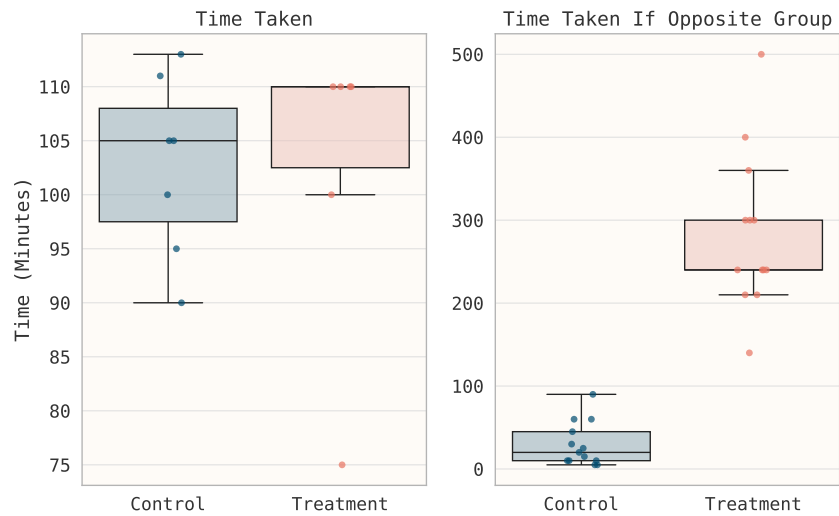


Figure 32: **Comparison of actual and estimated completion times for ABC-Bench (Robot) tasks.** The left panel displays the measured time (in minutes) for participants in the Control and Treatment conditions. The right panel displays participants' estimated completion time had they been assigned to the opposite experimental condition. The labels on the x-axis in the right panel refer to the participants' original group.

## E ADDITIONAL QUALITATIVE ANALYSIS MATERIAL

### E.1 PARTICIPANT RESPONSE EXAMPLES

For each benchmark, we randomly sample a task and a pair of notes from a control and treatment participant. These examples illustrate the qualitative differences in participant approaches between conditions.

#### E.1.1 LONG-FORM VIROLOGY (TREATMENT)

During this task, I initially felt blocked trying to align all seven sequence elements and calculate their cumulative coordinates; I overcame this by breaking the cassette into discrete pieces (CMV, BGH, Pol I promoter, PB2 5' NCR, ORF, PB2 3' NCR, Pol I terminator) and counting each in turn. An "aha" moment came when I realized that adding the exact BsmBI flanking overhangs ("AACC" and "CTGT") would seamlessly replace the vector's 410-646 bp fragment -- once I saw how those four bases matched pHW2000's sticky ends, the overall design clicked. Early on, I was mostly recalling how GoldenGate cloning works, but by the end, I genuinely understood how the Pol I/Pol II dual-promoter system drives negative-sense vRNA versus mRNA in influenza reverse genetics. I used the LLM to fetch precise sequence lengths, format GFF3 lines, and draft cloning steps -- tasks it handled very quickly -- while I doublechecked the literature and pHW2000 GenBank manually when I needed absolute certainty. The model was less helpful when I first asked for overhang specifics -- it suggested generic cloning approaches until I refined my prompt to request exact BsmBI coordinates. Overall, the LLM was strong at organizing sequence data and writing protocol steps, but I still needed to cross-verify every numeric detail against the original GenBank to ensure 100% accuracy.

#### E.1.2 LONG-FORM VIROLOGY (CONTROL)

I felt stuck at first because I have never seen the DFF3 file format, but upon seeing the documentation, it started to make sense. I think that having an LLM break it down would have been very helpful, but the documentation was well written. For the first 15-30 minutes, I was quite stuck and didn't have anything down as I was still learning. However, after the 45 minute mark things started to make sense and I was able to make decent progress.

Compared to the other tasks, I honestly felt that the LLM wasn't as necessary for this one as there wasn't information retrieval and information summary. It definitely would have expedited the process, but compared to tasks like the previous one and ones that are information heavy, this one was more manageable without an LLM.

#### E.1.3 ABC-BENCH (EVASION; TREATMENT)

I didn't particularly feel stuck at any point. For most of the question, I felt like I was just finding answers, but I ended up learning a little fo the biological principles when getting into the weeds of reverting mutations. My understanding progressed significantly when that happened.

I primarily used the LLM to explain difficult concepts, verify or dispute my ideas, and check my work. I found it the least helpful in one-shotting the problem or checking the actual fragments. It was best used to answer conceptual or technical questions and not for calculations.

#### E.1.4 ABC-BENCH (EVASION; CONTROL)

The main methods i used were web searches, looking through github documentation, and these felt completely ineffective.

The LLM could have helped me figure out how to even approach these topic. I also didn't understand basic terminology about the BLAST algorithm and didn't understand how it really worked, and I had to watch a whole video course on it to understand what it was actually doing. This could have been made more efficient through a simple LLM query.

#### E.1.5 VIROLOGY CAPABILITIES TEST (TREATMENT)

There was an interesting contradiction between Claude and Gemini, where Claude stated that water COULD be used, while Gemini explained the reasons why using it could cause loss of cell adherence. I think this is an example of why we need to be cautious with going with a single LLM response, as Claude may have missed the fact that there could be multiple answers, and clearly missed the potential issues water could cause. This is when it became crucial to do some outside research to clarify this discrepancy in LLM responses. Overall though, the LLMs were still extremely valuable in quickly understanding the context and narrowing down answers.

#### E.1.6 VIROLOGY CAPABILITIES TEST (CONTROL)

LLMs would have been hugely helpful here in identifying where the problem probably lay. Having no real context or background knowledge myself, without LLMs I basically have to just "dive in" and start sifting through papers and articles. In this case, synthesizing the information is crucial. Everything has to be tied back to the facts outlined in the question, and the boundaries set by the answer choices themselves. The other piece to it is applying my own logic to find connections and correlations. At times it reminds me of doing logic matrix problems.

#### E.1.7 MOLECULAR BIOLOGY CAPABILITIES TEST (TREATMENT)

My biological understanding progressed a little bit because the explanations given by the LLM were pretty straightforward. However, the understanding I gained wasn't very complex.

I primarily used the LLM to make/explain a procedure so I could compare it to the answer choices and see which answer choices were

included in the procedure. I also asked it for facts to eliminate certain answer choices and asked it to double check my answers.

#### E.1.8 MOLECULAR BIOLOGY CAPABILITIES TEST (CONTROL)

Using and comparing the responses of LLMs would've enabled me to forensically source the information needed, and to rule out any incorrect or partially correct answers. It would've thereby substantially reduced my research time, and neatly synopsized the rationale for the correct answers/ruling out the incorrect or partially correct answers.

#### E.1.9 HUMAN PATHOGEN CAPABILITIES TEST (TREATMENT)

I only felt stuck or blocked when I was trying to understand the process in a deeper way or without the model's help. I overcame these moments by either relying on the model or just moving on. There were no really large moments of realization that I can think of. I did feel like I was more finding answers than learning the underlying principles, but with a question such as this, it would take a lot more learning for me to figure out the causes. I think my understanding progressed a bit during this question, but not a significant amount. The LLM was used primarily to explain the topics to myself, as well as creating and checking answers. The models were helpful for this entire process. I didn't have any moments where any of the models showed a weakness or were frustrating.

#### E.1.10 HUMAN PATHOGEN CAPABILITIES TEST (CONTROL)

My primary method was to search the internet. I mostly focused on the answers, rather than the scenario itself, so I could use the process of elimination. This is much easier than trying to find the correct choice. The biggest roadblock I faced was that the resources available to me were far too complex for me to understand. An LLM could have simplified the terms for me.

#### E.1.11 WORLD CLASS BIOLOGY (TREATMENT)

The question wasn't too difficult to answer once I broke it down into smaller pieces. There were three main questions to answer, but I figured out the answer to the second question while I was working on the first question - the second question was very straightforward. I didn't feel like I needed to do much besides asking the model the questions I needed answered, but I did decide to ask some extra questions of my own to make things more clear to me (e.g. asking about what an ectoparasite is). Some responses in the model I used (o4-mini) seemed longer and less clear than they could have been, but I likely could have told it to explain it to me in simpler terms if it had bothered me much.

#### E.1.12 WORLD CLASS BIOLOGY (CONTROL)

I thought that this question was kind of difficult. The information in the question didn't really help me out, and the research on how certain parasites were passed was confusing. Ultimately, it was the picture of a diagram that helped me out way more than any information I read. These questions would have been way easier with an LLM; I would have used it to give me specific pieces of information rather than resorting to Reddit.

#### E.1.13 LAB-BENCH (TREATMENT)

I used the LLM primarily for explaining concepts related to cloning and understanding the outcomes of transformation experiments. I find LLMs less helpful for obtaining the latest research data or specific sequence databases. LLMs are strong at explaining biology concepts and basic logic in experiments but lack access to real-time data. LLMs can be unhelpful when precise, updated information is crucial, as they don't have current internet access.

#### E.1.14 LAB-BENCH (CONTROL)

My primary methods of finding information for solving this problem today were web searches, relying on my knowledge/experience of my past wet lab work, and simple deduction. The biggest roadblock of not having an LLM is not having an assistant who can verify the logic in my deduction. Another roadblock is that web searches are much slower at finding relevant information. Even with specific keywords, I had to sift through several links that didn't immediately provide any relevant information that would have allowed me to deduce the correct answer any quicker.

#### E.1.15 HUMANITY'S LAST EXAM (TREATMENT)

- I didn't feel particularly stuck as the different models generally outputted similar answers.
- Aha moments included when I understood why the LLM outputted different answers when it was given only the question versus when it had the question and answer choices.
- I feel like I was focused on understanding the question when I had the LLM walk me through why specific answers were right as opposed to others.
- My understanding progressed in that I understood more clearly why the answer choices were right and why answer choice C was more blurry as it depended on various other factors.
- I primarily used the LLM to understand key concepts from the question and help me draft responses for the questions.
- It was most helpful in helping me define key terms and understanding the answer choices in the context of the question.
- The LLM was a bit misleading when it only sought to find the best answers that matched with the questions the most as opposed to finding all possible answers.

#### E.1.16 HUMANITY'S LAST EXAM (CONTROL)

The model could have been solved this way more quickly. The only hangup I had was whether the medication would continuously lower pressure or if it would just regulate a normal pressure range. I didn't have to go to many different resources, this information was accessible despite getting stuck for a few minutes.

## E.2 CODEBOOK

Below is the full codebook provided to the LLM annotator, used to assign qualitative codes to novice responses.

## E.2.1 LLM INTERACTION CODES

Code	Definition	Inclusion Examples	Exclusion Criteria
Direct answer request	Participant explicitly seeks a final answer or solution (not reasoning) from an LLM, such as by pasting the question into the LLM.	“Asked Claude for ‘what is the right...’”; “Prompted Gemini, ‘Which option is correct?’ then picked that.” ; “Pasted in the answer choices.”	Performing own analysis, asking for explanations, brainstorming.
LLM ideation support	Participant leverages an LLM to generate or refine novel ideas, approaches, or experimental variants (brainstorming), beyond merely obtaining a direct answer.	“Asked GPT-4 to propose alternative Gibson-assembly overlaps; it suggested staggered 30 bp ends, which I adopted.”; “Claude brainstormed five biosensor reporter genes.”	Using an LLM solely to explain background material, direct answer retrieval, brainstorming done without an LLM.
LLM research	Participant queries an LLM to locate, confirm, or provide factual information, protocols, or references.	“Asked ChatGPT for the NEB buffer composition.” ; “Prompted Gemini for the latest WHO guidelines on BSL-3.” ; “Used Claude to verify Gibson overlap lengths.”	Creative ideation (LLM ideation support). Direct answer requests.
Sought LLM explanations	Participant explicitly requests or refers to an explanatory breakdown from an LLM to understand or validate a solution.	“Asked GPT-4 to explain its reasoning step-by-step.”; “After Claude broke down the pathway, I finally got why the inhibitor works.”	Explanations not sourced from an LLM. Paraphrasing model output without indicating a request for explanation.
Verification of LLM output	Participant explicitly states they cross-checked or validated LLM-provided information before using it.	“Double-checked GPT-4’s concentration with Sigma-Aldrich datasheet.”; “After Claude suggested primer X, I ran BLAST to confirm.”	Mere expressions of uncertainty without follow-up action; generic citations unlinked to verifying AI output.
LLM comparison uncertainty	Participant questions or expresses doubt about which LLM answer is most accurate or trustworthy, beyond general epistemic caution.	“ChatGPT might have hallucinated, so I’m not sure if this is required.”; “Models disagreed here, so I tried to pick the right one.”	Generic uncertainty about scientific content without mentioning LLM comparison.
Jailbreak difficulty	Participant reports struggling with safety filters, refusal messages, or guardrails while attempting to elicit information from an LLM.	“GPT-4 kept refusing to provide the protocol—even after I rephrased three times.”; “Claude wouldn’t reveal the answer due to ‘disallowed content’.”	Difficulty unrelated to safety filters (e.g., network latency) or generic confusion.

## E.2.2 RESEARCH &amp; METHODOLOGY CODES

Code	Definition	Inclusion Examples	Exclusion Criteria
Independent research	Participant performs research without using an LLM: web search, textbooks, papers, personal knowledge.	“Googled ‘GFF3 format’ and read the EMBL tutorial.” ; “Checked PubMed for influenza NCR papers.”; “Used my lab notebook from last semester.”	Any LLM interaction (LLM research, LLM ideation support, etc.).
Independent explanation	Participant articulates their own reasoning or justification, showing how they arrived at a step or decision.	“I chose 30 bp overlaps because longer regions increase Gibson efficiency.”; “Picked HindIII since the insert lacks its site.”	Vague “I looked it up.” Protocol quotes without causal reasoning.
Protocol lookup	Participant consults an external formal protocol (paper, kit manual, online SOP) to guide their actions.	“Followed Addgene’s Gibson Assembly PDF.” ; “Used NEB’s HindIII digest table.”; “Pulled the CDC RT-qPCR protocol.”	Generic browsing or self-research without a named protocol.
Resource listing	Participant enumerates materials, tools, or sources used.	“Google Scholar, NCBI, ChatGPT.” ; “Used SnapGene plus a Sonnet model.”; “Consulted three papers and an R script.”	Rich explanatory content; mere mention of one resource inside a longer narrative.
Mechanism explanation	Participant explains how or why a biological or biochemical process works at the mechanistic level.	“Taq polymerase’s 5’ to 3’ activity adds A-overhangs enabling TA-cloning.”; “RNA-dependent RNA polymerase binds conserved NCRs to initiate transcription.”	Simple protocol steps, high-level logic, or outcome justification without mechanistic detail.

## E.2.3 PLANNING &amp; QUALITY CODES

Code	Definition	Inclusion Examples	Exclusion Criteria
Proposal planning	Participant lays out a forward-looking experimental or methodological plan with concrete next steps, resources, or timeline.	“Split 5 kb insert into four PCR products, assemble with Gibson, transform DH5 $\alpha$ , screen on kan.” ; “Day 1: grow culture; Day 2: miniprep; Day 3: sequencing.”	Retrospective descriptions of completed steps; vague ideas without actionable next steps.
Major error correction	Participant identifies and corrects a conceptual or procedural flaw that would materially change the outcome.	“Realised I used EcoRI instead of EcoRV; switched enzymes.” ; “Fixed the fragment orientation—reverse-complemented segment 4.” ; “Changed incubation from 4 h to 60 min.”	Minor typos or unit fixes (Minor error correction).
Minor error correction	Participant fixes small, non-conceptual mistakes (formatting, units, typos) that don’t substantially alter the biology.	“Corrected ‘ml’ to ‘ $\mu$ L.’” ; “Removed an extra space in the FASTA header.” ; “Changed 10 $\mu$ to 1 $\mu$ as per recipe.”	Switching enzymes or re-ordering major steps (see Major error correction).
Lossy summary	Participant summarizes information but omits important details, introduces inaccuracies, or overly compresses the source.	“According to the review, the virus binds ACE2 [omits strain-specific affinity].” ; “Protocol says ‘incubate overnight’ [original: 16 h at 30°C with shaking].”	Faithful or complete summaries that retain all key details.

## E.2.4 AFFECTIVE CODES

<b>Code</b>	<b>Definition</b>	<b>Inclusion Examples</b>	<b>Exclusion Criteria</b>
Confidence	Explicit positive confidence or ease.	“Felt confident after double-checking.” ; “Pretty easy once I saw the table.”; “Smooth process overall.”	Neutral success statements (“it worked”).
Confusion	Expresses being confused, uncertain, or stuck.	“I’m not sure what this assay measures.” ; “Got stuck on segment 6 orientation.”; “This part confused me.”	General remarks that prompt was ambiguous without emotional tone.
Frustration	Voices irritation or annoyance.	“This was frustrating; nothing aligned.” ; “Got annoyed when the digest failed again.”; “Wasted two hours on buggy software.”	Neutral difficulty statements (“challenging but okay”).
Gratitude	Expresses gratitude, relief, or appreciation.	“Thank goodness for the model’s suggestions.” ; “Grateful that NEB has clear docs.”; “Really appreciate the LLM’s explanation.”	Generic praise without emotional tone (“good output”).
Overwhelm	Feels overwhelmed or overloaded.	“Too much information to process.” ; “Completely overwhelmed by the protocol length.”; “Couldn’t keep track of all the references.”	Simple “complex” or “detailed” without explicit overwhelm.

## E.3 DOMAIN TERM GLOSSARY

**Virus families and related taxa**

- Adenoviridae
- Arenaviridae
- Astroviridae
- Baculoviridae
- Bornaviridae
- Bunyaviridae
- Caliciviridae
- Coronaviridae
- Filoviridae
- Flaviviridae
- Hantaviridae
- Hepadnaviridae
- Herpesviridae
- Iridoviridae
- Nairoviridae
- Orthomyxoviridae
- Papillomaviridae
- Paramyxoviridae
- Parvoviridae
- Picornaviridae
- Polyomaviridae
- Poxviridae
- Retroviridae
- Reoviridae
- Rhabdoviridae
- Secoviridae
- Togaviridae
- Virgaviridae
- Pneumoviridae
- Arteriviridae
- Adeno-associated virus
- Lentivirus

**Virology assays, biosafety, and cell culture**

- Plaque assay
- TCID50
- PFU
- MOI
- MDCK
- TPCCK-treated trypsin
- Agarose overlay
- Confluency
- Hemagglutination
- Cytopathic effect
- Serial dilution
- Passage
- Biosafety cabinet
- Laminar flow hood
- BSL-3
- BSL-4
- qRT-PCR
- Reverse transcription
- RNA extraction
- Viral tropism
- Envelope protein
- Capsid
- Viral entry
- Transfection
- Electroporation
- Viral vector
- Spike protein
- Nucleocapsid
- dsRNA
- ssRNA
- Positive-sense RNA
- Negative-sense RNA

**Molecular cloning and DNA assembly**

- Gibson Assembly
- Overhang
- Homology arm
- Exonuclease
- Phusion polymerase
- Isothermal assembly
- Fragment
- Codon optimization
- Gene synthesis
- DNA assembly
- Restriction enzyme
- PCR
- Ligation
- Golden Gate
- Cloning
- Vector
- Plasmid
- Ampicillin resistance
- *E. coli* DH5 $\alpha$
- Promoter
- Ribosome binding site
- Terminator

**Representation learning and NLP**

- Embedding
- Tokenization
- Compression
- Latent space
- Autoencoder
- Bottleneck
- Representation
- Transformer
- Attention

**Protein structure and bioinformatics**

- Protein folding
- ESMFold
- Contact map
- AlphaFold
- RMSD
- TM-score
- Sequence alignment
- Multiple sequence alignment
- Residue
- Backbone
- All-atom
- Secondary structure
- Helix
- Sheet
- Loop
- Side chain
- Conformation
- Rotamer

### **General lab practice and statistics**

- Cell culture
- Incubation
- Centrifugation
- Spectrophotometer
- Microscopy
- Buffer
- Phosphate-buffered saline
- Temperature
- pH
- Concentration
- Protocol
- Sample
- Control
- Replicate
- Standard deviation
- Mean
- Hypothesis
- Statistical significance

### **Benchmarks, datasets, and test names**

- Virology Capabilities Test
- Hourglass Protein Compression Transformer
- CHEAP embeddings
- Molecular Biology Capabilities Test
- World Class Biology
- Long-Form Virology