

WETBENCH: LLM-BASED SIMULATION ENVIRONMENT TO EVALUATE WET-LAB EXPERIMENT PLANNING AND DESIGN

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce WetBench, an LLM-based simulation environment for scalably evaluating AI systems’ ability to design and plan wet-lab experiments. Traditional evaluation approaches are limited by the expense and safety concerns of executing AI-generated experiments in physical laboratories. To address this, we developed a simulation environment that uses LLMs as state transition models to simulate experimental outcomes and as evidence classifiers to evaluate whether experiments provide sufficient information to achieve stated goals. WetBench includes 18 expert-curated experimental configurations spanning cell biology, neuroscience, microbiology, and analytical chemistry, each validated as solvable within the environment’s constraints. We evaluated the fidelity of our LLM-based simulation through expert ratings, finding that state transitions were judged as highly plausible (90% plausibility) by human expert raters. Evidence classification showed substantial agreement between LLM classifiers and human experts (72-82% agreement), on par with inter-human baseline agreement (75%). Using this environment, we benchmarked frontier language models on experimental design and planning capabilities. GPT-5 demonstrated superior performance with a 44.4% pass@1 rate that increased to 72.2% at pass@5, substantially outperforming other models, including Gemini 2.5 Flash (50.0% pass@5), Qwen 3 (41.2% pass@5), and Claude Sonnet 4 (27.8% pass@5). We open-source WetBench as a Python gymnasium environment to support further development of AI systems for autonomous scientific experimentation. ¹

1 INTRODUCTION

Artificial intelligence has significant potential to accelerate research in biology and bioengineering. While recent work has demonstrated substantial progress in using AI systems for hypothesis generation (Gottweis et al.; Baek et al.) and data analysis (Aygün et al.), a critical component of the discovery process in biology remains underexplored: experimental planning and design. Historically, large language models (LLMs) have been criticized for weak planning capabilities: they often generate plans that are not executable, lack the ability to self-verify or refine them, and struggle to manage complex constraints (Kambhampati et al.; Xie et al.; Vyas et al.; Chang et al.). Additionally, while there is work suggesting reasonable human-LLM agreement in hypothesis evaluation (Ghareeb et al.), this has not been established in the context of evaluating experimental design. Progress in this area is critical, as we may soon be bottlenecked not by the number of scientific hypotheses AI systems can generate, but by our ability to design and execute experiments to evaluate them (Reddy & Shojae; Zhang et al.).

While successful experimental planning and design requires many features, we focus on two in particular:

- **Feasibility** A feasible experiment consists of a set of steps that can be performed in the lab, given constraints on materials, equipment, experimental protocols, and realistic expectations.

¹We will provide the link to the GitHub after de-anonymization

- **Informativeness** An informative experiment provides the information required to answer the intended question through well-designed controls and conditions.

Evaluating the ability of AI systems to develop feasible and informative experiments in the physical world is challenging for multiple reasons. Having human experts execute every LLM-generated experiment in the lab is prohibitively expensive, slow, and could confound AI experimental ability with human expertise. While many groups have used liquid handlers as a platform to explore autonomous experimental design, liquid handlers significantly constrain the space of potential experiments. On the other hand, providing AI agents with direct, unsupervised access to large, diverse experimental facilities—such as self-driving labs (Qiu et al.) or cloud laboratories (Boiko et al.)—raises safety concerns and incurs significant resource and time costs (Sandbrink). Each approach falls short of providing the scale and diversity necessary to meaningfully evaluate and improve experimental planning capabilities in AI systems.

Considering this, we sought to determine whether LLM-based experimental simulation environments have the potential to serve as surrogates for experiments. Through strong performance on challenging biological reasoning benchmarks like GPQA (Rein et al.) and Lab-Bench (Laurent et al.), there is evidence that frontier reasoning LLMs have internalized large amounts of biological and experimental knowledge. If we can use LLMs to realistically simulate experimental procedures and judge when experiments are sufficiently informative, we could provide a sandbox for evaluating LLM experimental design ability and an environment to train AI agents in a way that is scalable, highly diverse, and safe (such as in Liu et al.). Critically, for the system to be useful, we do not need or expect LLMs to predict completely novel experimental outcomes. Experimental design in practice involves executing well-established protocols to investigate questions where outcomes remain uncertain. For example, to test a new biosensor, a researcher generally plans by modifying experimental protocols they are confident in to design the experiment with appropriate controls. While LLMs should not be able to determine if the biosensor would work, they should be able to design feasible experiments using reasonable protocols and informative controls to gather the information required to answer the question.

To explore this idea, we developed WetBench, an LLM-based wet-lab simulation environment. AI agents begin with an experimental goal, initial materials, and available actions (e.g., Combine, Incubate, UseMicroscope). At each step, agents take actions on materials and submit them to a state transition model, which simulates resulting material property changes and physical observations. After accumulating a history of actions and observations, agents submit their evidence to a classifier that determines whether the cumulative results sufficiently justify the experimental goal. We manually curated a set of experimental configurations (i.e., an experimental goal and initial materials) to serve as a benchmark for the WetBench environment. After verifying that each experiment was solvable within the constraints of the environment, we had AI agents attempt to solve these problems and used the resulting state transitions and evidence submissions to evaluate the system’s fidelity. Our work extensively examines when these LLM-based state transitions and evidence interpretations align with expert judgment and identifies systematic failure modes.

Our work makes five key contributions: (1) we develop a simulation environment for wet-lab experiment execution; (2) generate an expert-verified experiment benchmark for evaluation; (3) evaluate the fidelity of LLM-based state transitions and evidence classification using expert ratings; (4) explore the current abilities of frontier and open-source models on experimental design within the environment; and (5) open source the project to support further development of both the environment and future AI agents.

2 RELATED WORK

A complementary line of work focuses on translating human-composed protocols into structured instructions. BioPlanner (O’Donoghue et al.) converts natural language protocols into pseudocode for laboratory automation, utilizing a teacher-student framework. In their approach, a “teacher” LLM model converts existing protocols from the literature into pseudocode for execution in a robotic lab. A “student” model is then evaluated on its ability to reconstruct this pseudocode when given only the protocol description and admissible pseudofunctions. The student’s pseudocode is then evaluated based on next-step prediction, full protocol generation, and pseudofunction retrieval, with success measured against the teacher model’s pseudocode. While this approach enables automatic

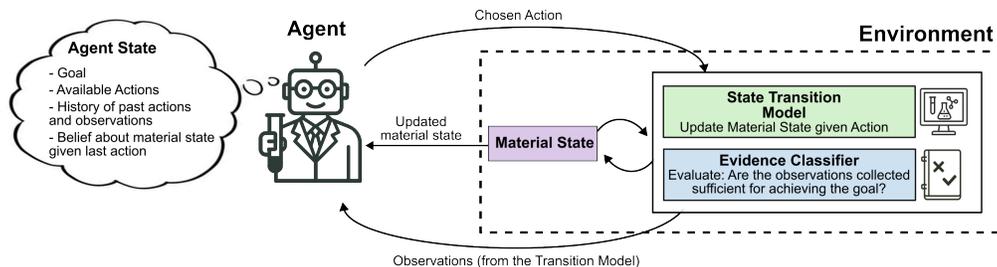


Figure 1: **Overview of the WetBench system.** The AI agent maintains goals, prior history, and knowledge of the available materials and actions. At each step, the agent chooses an action to make progress towards its goal. If the action is a *transformation*, *measurement*, or *simulation* action, it is submitted to the state transition model that determines the result of new materials and observations. If the agent chooses "Submit" then all prior actions and observations are sent to the evidence classifier to determine if the goal was achieved.

evaluation of protocol conversion accuracy, it addresses the translation of known experimental procedures rather than the design of novel experiments to achieve informational goals. Additionally, their framework does not model experimental outcomes or assess whether proposed actions would actually succeed in achieving the scientific objectives. In contrast, our work addresses open-ended experimental planning, where models must determine which experiments to perform based on scientific objectives, with success defined by the sufficiency of evidence for target interpretations rather than adherence to predetermined procedural sequences.

Coscientist (Boiko et al.) demonstrated autonomous experimental planning, using LLMs with tool access, including web search, code execution, documentation retrieval, and direct control of robotic lab equipment through APIs. The execution through robotic lab systems (such as Opentrons and Emerald Cloud Labs) defines a fixed action space that constrains the plans of the LLM coscientist. However, in this system, the goals given to the coscientist are always assumed to be feasible given the resources and action space, ignoring the important roles of feasibility and evidence sufficiency evaluation in planning.

3 WETBENCH SYSTEM OVERVIEW

The WetBench system consists of an environment and an agent. The environment is comprised of the material state, associated observations, a set of experimental actions, a state transition model that simulates these actions, and an evidence classifier that determines whether the series of materials, actions, and observations constitutes sufficient evidence for the experimental goal.

The agent operates with an experimental goal and knowledge of available actions and materials in the environment. The agent's task is to design and execute an experiment in the environment that is sufficiently informative to achieve the experimental goal. These components interact in a closed-loop fashion: the agent attempts to execute actions in the environment; the state transition model either rejects the action (due to it being infeasible or outside of the set of available actions) or simulates the effect of the action by modifying the relevant materials, producing physical observations, and stepping in time. After receiving the observation and materials changes, the agent chooses experimental actions until it believes it has reached its goal. When the agent believes its actions and observations are sufficient for the goal, it submits the results to the evidence classifier, which accepts or rejects them based on whether they are sufficiently informative.

3.1 AGENT

The agent is an LLM-based system prompted to achieve specific experimental goals using a defined set of available materials and actions. The agent receives an experimental objective, an initial inventory of materials, and knowledge of the available action space. For instance, the goal might be to determine calcium response dynamics in cortical neurons following KCl depolarization, starting with materials like specific plasmids, stock solutions, and neuronal cultures.

The agent’s decision-making process relies entirely on the LLM’s reasoning capabilities and scientific knowledge, without access to external tools or databases beyond what is explicitly provided in the simulation environment. This design ensures that experimental planning performance reflects the model’s intrinsic understanding of laboratory procedures and scientific methodology.

3.2 MATERIAL STATE AND ACTIONS

Material state The material state defines what the agent has access to conduct its experiments. Typical materials include biological samples (e.g., cell lines, primary neuron cultures, etc.), reagents (e.g., buffers, growth media, etc), containers (e.g., Eppendorf tubes, 50 mL conical centrifuge tubes, etc.), and chemicals (e.g., hydrochloric acid, sodium hydroxide, ethanol, etc.). Each material is represented using a structured dictionary format with a name, physical/chemical/properties properties (e.g., concentration, container specifications, cell type, etc.), and a unique barcode identifier (Figure 2, Top Left). Various addition properties are included to enforce realism in the environment. Each material has a defined environmental condition in which it exists (e.g., held at room temperature, in a 4°C refrigerator, or -80°C freezer, etc). Finally, each material is labeled as static or dynamic. Dynamic materials are materials that dynamically change over time; for example, cells that continue to grow or chemical reactions that require a certain amount of time to complete. The age of the material is tracked through the "created_at" and "last_modified" properties. As time passes in the simulation environment, the material state changes based on the environmental conditions it is in and its properties.

Actions At each step, the agent chooses an action to apply to a subset of materials to make progress towards the goal. The set of actions was designed to cover the set of typical actions done during wet-lab experiments (detailed in Appendix C). Some actions like "Transfer," "Incubate," or "Wash" are *transformation* actions as they primarily change the physical properties of the materials. Other actions like "UseMicroscope" or "MeasurepH" are *measurement* actions because they primarily produce data and observations (they may or may not affect the physical properties of the material). Finally, there are *simulation* actions like "Wait" that allow the agent to interact with its simulation environment directly (e.g., stepping forward in time). Each action has required and optional parameters with fixed constraints on the parameter values (Figure 2, Bottom Left). This is to prevent the agent from inventing action parameters to "hack" the system into producing any output it wants.²

3.3 STATE TRANSITION MODEL

When the agent chooses an action, parameters, and target materials, this information is submitted to the State Transition Model (detailed in Appendix C.2). The State Transition Model is an LLM prompted to decide if the action should be rejected and, if not, to simulate plausible outcomes of the action. The rejection reasons are not limited to, but include:

- Missing critical materials or reagents
- Insufficient amounts or values (e.g., the new container intended to hold the sample is too small, the measurement device cannot reliably record a signal in liquid below a certain volume)
- Incorrect material state (e.g., frozen samples needing to be in a liquid state, adherent cells needing to be in suspension)
- The described action is not valid considering the set of available actions and parameter constraints.

When simulating plausible outcomes, the state transition model is prompted with a detailed description of how the action would be conducted in the physical world and instructed to consider:

- Chemical reactions and physical processes
- Conservation laws and material balance (i.e., what goes in must come out)

²"VisualInspection" is the only action that has a free-string parameter (it is constrained to be a regular expression that defines a single sentence that ends with a question mark). However, the state transition model is prompted to reject any inquiries that cannot be ascertained by looking at the sample.

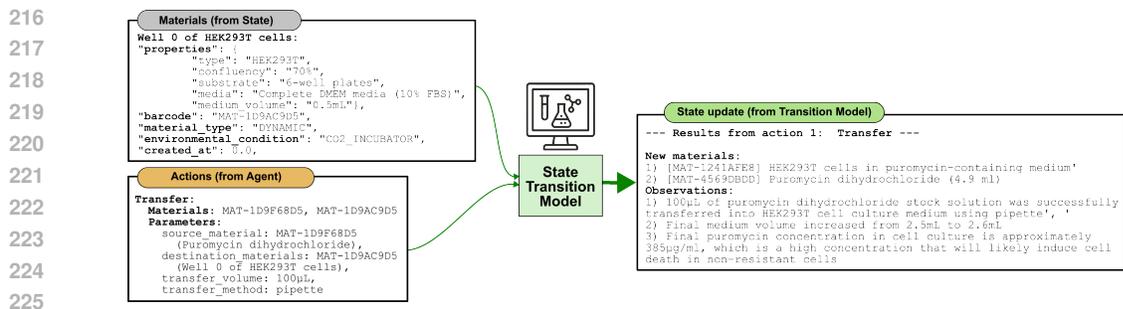


Figure 2: **State transition model workflow.** Upper left: Example of individual material in material state with a set of physical properties, a unique barcode, and environmental information. Another material (puromycin) is omitted due to space constraints. Lower left: Example of action. The target materials refer to the material that goes into the action (e.g., HEK293T Cells and puromycin). The action has associated parameters including the source and destination materials, the transfer volume, and the transfer method. Right side: The outcome of a state transition model for the action. Two new materials are generated: HEK293T cells with puromycin and residual puromycin. Additionally, the state transition model generates a set of physical observations corresponding to the action and the new materials that have been introduced.

- Realistic experimental outcomes
- Time-related processes (growth, decay, transformation, etc.)
- Observational limits (what could be observed by a human or AI with human-like sensory capabilities.)

The state transition model outputs a structure dictionary that describes the resulting material state and a set of associated observations from the experimental action. It is through these constraints on the state transition model that we work to force the agent into producing **feasible** experiments.

3.4 EVIDENCE CLASSIFIER

Within the environment, an experiment can end for one of three reasons:

- The agent believes it has sufficient evidence for its goal and executes the "Submit" action, which submits the history of observation to the evidence classifier.
- The agent does not believe it can accomplish its goal and executes the "Quit" action
- The agent attempts the maximum number of allowed actions (default 100 actions), or three actions are rejected in a row

The evidence classifier evaluates whether cumulative experimental evidence supports the stated goal by receiving a comprehensive experimental history, including all actions, materials, timing, and collected observations, organized chronologically. This evidence history is evaluated against the experimental goal. The classification prompt instructs the model to consider scientific plausibility, logical consistency, appropriate inference levels, alternative explanations, evidence strength, and confounding factors in its reasoning (detailed in Appendix C.3).

For example, in an experiment aimed at determining calcium response dynamics in cortical neurons following KCl depolarization, sufficient evidence might include: baseline calcium measurements, controlled KCl application with appropriate concentrations, time-series calcium imaging data showing clear signal changes, and negative controls without KCl treatment. Insufficient evidence might consist of only a single calcium measurement without controls, unclear timing of KCl application, or missing baseline measurements that prevent interpretation of any observed changes. Critically, it is through the evidence classifier that we enforce the **informativeness** requirement of successful experimental design. Just as the state transition model constrains agents to produce feasible protocols by rejecting impossible actions, the evidence classifier ensures that agents must gather sufficient evidence—with appropriate controls and logical inference—to meaningfully address their experimental goals, rather than simply executing a series of plausible-sounding steps.

4 METHODS

4.1 EXPERIMENT CONFIGURATIONS

For the benchmark, we generated 18 experimental configurations spanning cell and molecular biology, neuroscience, microbiology, and analytical chemistry. Each experimental configuration is defined by a set of materials and a scientific goal. We curated the experimental configurations based on experiments that the authors have run in the wet-lab or based on the literature. The full list of experiments is in the (detailed in Appendix B).

4.2 MODEL CONFIGURATION

We evaluated multiple models as experimental design agents:

- **Proprietary reasoning models:** Claude 4 Sonnet (extended thinking), GPT-5, Gemini 2.5 Flash
- **Open-source reasoning models:** Qwen 3 (235B), DeepSeek R1
- **Non-reasoning model:** GPT-4o

For the WetBench environment components, we used Claude 4 Sonnet as the state transition model and GPT-5 as the evidence classifier. Complete model parameters are provided in the Appendix A.

4.3 EXPERT RATINGS

For each evaluation task, we selected multiple Ph.D. or postdoc-level researchers in biology, bio-engineering, or neuroscience. We instructed each evaluator on the specifics of the tasks and reviewed two to three examples prior to the ratings. Each rating generates a binary score (either plausible or not plausible for state transitions and sufficient or not sufficient for evidence classification), as well as a confidence measure from one (very low) to five (very high).

5 RESULTS

5.1 STATE TRANSITION MODEL VALIDATION

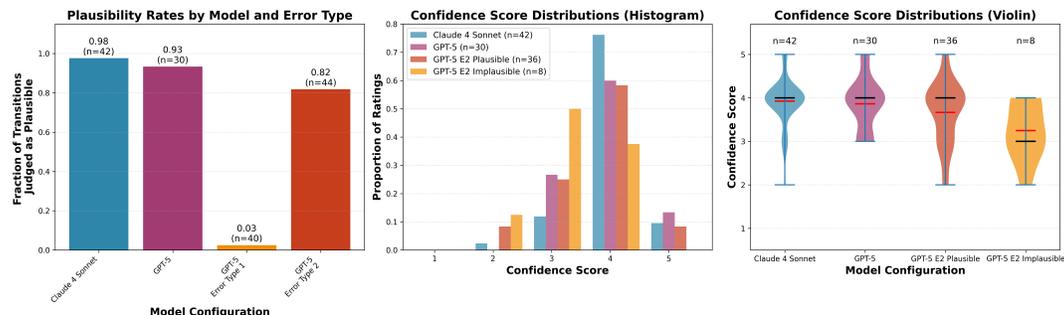


Figure 3: **Expert validations of state transition model.** Left Panel: Across ratings, we find overwhelming agreement that the state transition tests are plausible for Claude 4 Sonnet and GPT-5. Experts rate GPT-5 error type 1 transitions extremely low and rate GPT-5 error type 2 at lower but similar plausibility rates to the unperturbed models. Center Panel: Histogram distribution of the Claude 4 Sonnet confidence ratings, GPT-5 confidence ratings, and GPT-5 error type 2 confidence ratings split into ratings that were deemed plausible and implausible. Right Panel: A violin plot comparing Claude Sonnet, GPT-5, and GPT-5 error type 2 split into plausible and implausible.

To understand the performance of different LLMs as state transition models, we collected a large set of (input materials, action) pairs from multiple simulation runs and used Claude 4 Sonnet and GPT-5 to generate state transitions. We then collected plausibility ratings from evaluators, who were

Table 1: Inter-rater agreement metrics between human-model and human-human

MODEL	AVG KAPPA	AVG AGREEMENT
Claude 4 Sonnet	0.618	81.9%
GPT-4o	0.583	80.4%
Gemini 2.5 Flash	0.472	74.7%
GPT-5	0.437	72.2%
Qwen 3	0.430	72.2%
Human	0.526	75.0%

prompted to assess whether the outcome materials and observations were plausible given the action and input materials.

LLM-simulated state transitions rated as highly plausible. Evaluators found the overwhelming majority of predicted transitions to be plausible (90% for Claude 4 Sonnet, 93% for GPT-5). Confidence was generally high (average of 3.90 ± 0.59), and anecdotal reports suggested that both models could produce reasonably plausible outcomes. Unsurprisingly, due to the high agreement on plausibility, inter-rater agreement was quite high (Cohen’s $\kappa = 0.811$).

Raters struggled to identify subtle errors in state transitions. To calibrate our results, we used GPT-5 to introduce intentional errors into state transitions using two error types: obvious errors (Error Type 1) and subtle but legitimate errors (Error Type 2) (see Appendix C.2). This approach allowed us to evaluate rater performance when transitions were intentionally incorrect across different error sensitivities. While raters easily identified Error Type 1 transitions (3% rated as sufficient), their judgments of Error Type 2 transition plausibility were similar to baseline transitions. In our post-hoc analysis, we found that Error Type 2 transitions contained very subtle errors, such as calculation errors in reagent concentration, incorrect magnification objectives for imaging, or adding frozen reagents before thawing. These errors were typically single mistakes embedded within heavily detailed and otherwise reasonable outcomes. Given that evaluators consistently reported that parsing state transitions was challenging due to the cognitive load, they appear reliably able to discern general plausibility but struggle with subtle technical errors.

Subtle errors are correlated with low confidence ratings Given the cognitive load of parsing dense state transition information, we examined whether type 2 errors correlated with rater confidence. We compared confidence score distributions between regular predictions and type 2 error predictions. There was a modest relationship between confidence and correctness (AUC = 0.609), though not significant enough to serve as a reliable predictor (Figure 3). This relationship held for both error types regardless of whether they were judged as plausible or implausible. However, transitions judged as implausible showed substantially lower mean and median confidence scores compared to all other conditions.

5.2 EVALUATING EVIDENCE CLASSIFICATION AGREEMENT

Evaluators agree with LLM and human evidence classifiers at similar rates. To validate LLM-based evidence classification, we compared agreement rates between human expert evaluators and different language models. Claude 4 Sonnet achieved the highest agreement with human evaluators at 81.9% ($\kappa = 0.618$), followed by GPT-4o at 80.4% ($\kappa = 0.583$). Both models exceeded the baseline inter-human agreement of 75.0% ($\kappa = 0.526$). The remaining models showed lower but still substantial agreement: Gemini 2.5 Flash at 74.7% ($\kappa = 0.472$), and both GPT-5 and Qwen 3 at 72.2% ($\kappa = 0.437$ and 0.430, respectively). The narrow range of agreement rates across models (72.2% to 81.9%) and the fact that top-performing models exceeded human baseline agreement suggest that LLMs can serve as reliable evidence classifiers in the simulation environment without requiring extensive human expert evaluation for each assessment.

Inter-model agreement reveals systematic differences in classification stringency. Overall agreement on evidence sufficiency among the five models was high, with agreement rates rang-

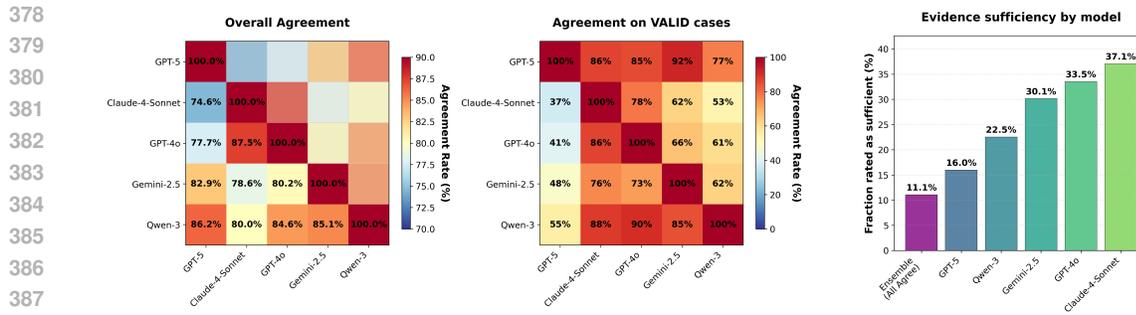


Figure 4: **Inter-model agreement rates.** Left: Agreement matrix between each model for both evidence judged sufficient and insufficient. Middle: Conditional agreement plot showing the rate at which the column model believes evidence is sufficiently valid, conditional on the row model believing that evidence is sufficiently valid. Right: The overall fraction of evidence deemed sufficient across different models and the ensemble. Note: DeepSeek R1 was excluded from this analysis due to persistent errors in its API.

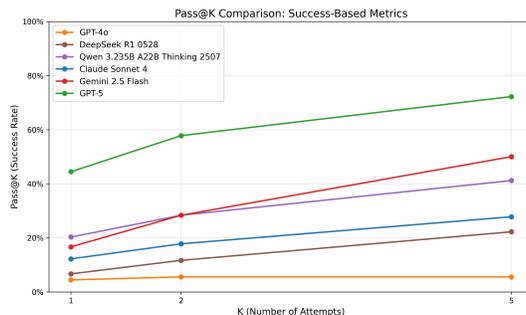


Figure 5: **pass@k rates across all models.** Increasing the number of attempts (K) leads to increased performance across all reasoning models, but does not improve GPT-4o’s performance.

ing from 74.6% to 87.5%. However, this is somewhat inflated by overwhelming agreement on cases in which there is insufficient evidence. When examining cases where models disagreed on sufficient experiments, clear patterns emerged in their classification tendencies (Figure 4). GPT-5 demonstrated the most conservative approach, classifying only 16% of evidence submissions as sufficient—the lowest rate among all models tested. In experiments that GPT-5 judged to be sufficient, all other models graded the evidence as sufficient at rates between 77% to 92% (Figure 4, Middle Panel). Conversely, Claude 4 Sonnet was the most permissive classifier, rating more than a third of experiments run as having produced sufficient evidence. When restricting analysis to cases where all models reached consensus on sufficiency, only 11.1% of total classifications were deemed adequate. This suggests significant heterogeneity in what experiments different LLM evidence classifiers believe are sufficiently informative.

5.3 FRONTIER MODEL PERFORMANCE IN THE WETBENCH ENVIRONMENT

Using the WetBench environment and benchmark, we evaluated how well frontier language models could accomplish experimental goals across the 18 experimental configurations. Informed by our validation results, we used Claude 4 Sonnet as the state transition model and GPT-5 as the evidence classifier for all agent evaluations. Although GPT-5 did not achieve the highest alignment with human raters, we utilized it to evaluate our metrics due to its overall evidential stringency.

GPT-5 demonstrates superior experimental design capabilities. GPT-5 achieved the strongest performance across all metrics, with a pass@1 rate of 44.4% that increased to 72.2% at pass@5 (Figure 5). This substantial performance represents successful experimental design and execution

432 on nearly three-quarters of the solvable benchmark problems when given five attempts. The model’s
433 high baseline success rate of 44.4% indicates consistent experimental planning abilities rather than
434 reliance on multiple attempts.

435
436 **Performance stratifies into distinct tiers across model families.** The remaining models showed
437 considerable performance variation, clustering into distinct capability tiers. Gemini 2.5 Flash and
438 Qwen 3 formed a second performance tier, achieving pass@5 rates of 50.0% and 41.2% respec-
439 tively, though both started from lower pass@1 baselines (16.7% and 20.3%). Claude Sonnet 4
440 and DeepSeek R1 demonstrated more limited capabilities, with pass@5 rates of 27.8% and 22.2%.
441 Notably, Claude Sonnet 4’s constrained performance may reflect its limited thinking budget (2048
442 tokens) relative to other reasoning models rather than fundamental planning limitations. GPT-4o
443 showed the weakest performance, achieving only a 5.6% pass@5 rate, suggesting that reasoning
444 capabilities may be particularly important for experimental design tasks.

445
446 **Qualitative analysis of AI experimental planning and design** To understand how successful
447 AI agents approach experimental design compared to human experts, we examined three represen-
448 tative cases where models achieved their experimental goals (see Appendix B for extended case
449 studies). Our qualitative analysis revealed distinct patterns in AI versus human experimental plan-
450 ning. For protocol-like tasks with minimal control requirements, AI agents converged with human
451 strategies—Gemini 2.5 Flash pursued nearly identical approaches to human experts when determin-
452 ing GFP insert orientation, differing only in minor procedural details. However, when experimen-
453 tal goals allowed multiple valid interpretations, strategies diverged substantially. In characterizing
454 acidic soil isolates, Qwen 3 employed repeated baseline pH measurements against each culture’s
455 own baseline, while human experts used uninoculated media controls—both representing reason-
456 able but distinct approaches to the same goal. Most notably, AI agents consistently missed domain-
457 specific shortcuts that experts readily recognized. When determining zebrafish genotypes, GPT-5
458 defaulted to DNA extraction workflows despite the availability of a simple visual inspection short-
459 cut that human experts employed, highlighting how domain expertise enables efficient experimental
460 design that general reasoning approaches may overlook.

461 6 CONCLUSION

462
463 In this work, we explored the potential of using LLMs to construct a simulation environment that
464 promotes high-quality experimental planning and design. Our validation studies demonstrated that
465 LLM-based state transition models produce outcomes that experts rate as highly plausible (90%
466 plausibility ratings), though we identified persistent challenges in evaluating subtle technical errors
467 within information-dense state transitions. For evidence classification, we found that LLM agree-
468 ment with human experts often matched or exceeded inter-human baselines (81.9% for Claude 4
469 Sonnet vs. 75.0% human-human agreement), suggesting strong potential for automated evidence
470 evaluation in experimental contexts. Using the WetBench environment to evaluate frontier models
471 revealed substantial performance differences, with GPT-5 achieving the strongest results at 72.2%
472 pass@5, demonstrating successful experimental design on nearly three-quarters of benchmark prob-
473 lems.

474 The WetBench environment enables rapid and scalable training of experimental design agents
475 through reinforcement learning, offering high-fidelity state transitions and reliable evidence clas-
476 sification. This virtual setting allows agents to iterate through thousands of experimental scenarios,
477 advancing planning strategies that would be impractical to develop in physical laboratories. As
478 hypothesis generation by AI outpaces our experimental capacity, WetBench offers a solution for
479 developing agents that can help close that gap and accelerate biological research.

480 REFERENCES

481
482 Eser Aygün, Anastasiya Belyaeva, Gheorghe Comanici, Marc Coram, Hao Cui, Jake Garrison, Re-
483 nee Johnston Anton Kast, Cory Y McLean, Peter Norgaard, Zahra Shamsi, David Smalling,
484 James Thompson, Subhashini Venugopalan, Brian P Williams, Chujun He, Sarah Martinson,
485 Martyna Plomecka, Lai Wei, Yuchen Zhou, Qian-Ze Zhu, Matthew Abraham, Erica Brand, Anna
Bulanova, Jeffrey A Cardille, Chris Co, Scott Ellsworth, Grace Joseph, Malcolm Kane, Ryan

- 486 Krueger, Johan Kartiwa, Dan Liebling, Jan-Matthis Lueckmann, Paul Raccuglia, Xuefei, Wang,
487 Katherine Chou, James Manyika, Yossi Matias, John C Platt, Lizzie Dorfman, Shibl Mourad, and
488 Michael P Brenner. An AI system to help scientists write expert-level empirical software.
489
- 490 Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. ResearchAgent: Iterative
491 research idea generation over scientific literature with large language models.
- 492 Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research
493 with large language models. 624:570–578.
494
- 495 Edward Y Chang, Arihant Choudhary, Parth Behani, and Yash Vardhan Pansari. Position: Limita-
496 tions of LLMs can be overcome by carefully designed multi-agent collaboration.
- 497 Ali Essam Ghareeb, Benjamin Chang, Ludovico Mitchener, Angela Yiu, Caralyn J Szostkiewicz,
498 Jon M Laurent, Muhammed T Razzak, Andrew D White, Michaela M Hinks, and Samuel G
499 Rodriques. Robin: A multi-agent system for automating scientific discovery.
- 500 Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom
501 Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici,
502 Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat,
503 Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan,
504 Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary
505 Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an
506 AI co-scientist.
- 507 Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant
508 Bhambri, Lucas Saldyt, and Anil Murthy. LLMs can't plan, but can help planning in LLM-
509 modulo frameworks.
- 510 Jon M Laurent, Joseph D Janizek, Michael Ruzo, Michaela M Hinks, Michael J Hammerling, Sid-
511 dharth Narayanan, Manvitha Ponnampati, Andrew D White, and Samuel G Rodriques. LAB-bench:
512 Measuring capabilities of language models for biology research.
513
- 514 Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone.
515 LLM+P: Empowering large language models with optimal planning proficiency.
516
- 517 Odhran O'Donoghue, Aleksandar Shtedritski, John Ginger, Ralph Abboud, Ali Essa Ghareeb, Justin
518 Booth, and Samuel G Rodriques. BioPlanner: Automatic evaluation of LLMs on protocol plan-
519 ning in biology.
- 520 Yibo Qiu, Zan Huang, Zhiyu Wang, Handi Liu, Yiling Qiao, Yifeng Hu, Shu'ang Sun, Hangke
521 Peng, Ronald X Xu, and Mingzhai Sun. BioMARS: A multi-agent robotic system for autonomous
522 biological experiments.
- 523 Chandan K Reddy and Parshin Shojaee. Towards scientific discovery with generative AI: Progress,
524 opportunities, and challenges.
525
- 526 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
527 Dirani, Julian Michael, and Samuel R Bowman. GPQA: A graduate-level google-proof Q&A
528 benchmark.
- 529 Jonas B Sandbrink. Artificial intelligence and biological misuse: Differentiating risks of language
530 models and biological design tools.
531
- 532 Kaustubh Vyas, Damien Graux, Sébastien Montella, Pavlos Vougiouklis, Ruofei Lai, Keshuang Li,
533 Yang Ren, and Jeff Z Pan. An extensive evaluation of PDDL capabilities in off-the-shelf LLMs.
- 534 Jian Xie, Kexun Zhang, Jiangjie Chen, Siyu Yuan, Kai Zhang, Yikai Zhang, Lei Li, and Yanghua
535 Xiao. Revealing the barriers of language agents in planning.
536
- 537 Yanbo Zhang, Sumeer A Khan, Adnan Mahmud, Huck Yang, Alexander Lavin, Michael Levin,
538 Jeremy Frey, Jared Dunnmon, James Evans, Alan Bundy, Saso Dzeroski, Jesper Tegner, and Hec-
539 tor Zenil. Exploring the role of large language models in the scientific method: from hypothesis
to discovery. 1:1–15.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

A APPENDIX

B CASE STUDIES

To understand how successful AI agents approach experimental design compared to human experts, we examined three representative cases where models achieved their experimental goals. These cases reveal both convergent and divergent strategies between AI and human experimental planning.

Case Study 1: Determine which of 20 colonies contain correctly oriented and in-frame GFP inserts without mutations

AI agents converge with human strategies on protocol-like tasks. Gemini 2.5 Flash pursued a nearly identical approach to the human expert. Both aliquoted multiple media tubes for storage, inoculated GFP colonies, incubated cultures over 24 hours, then extracted and sequenced DNA. The only notable difference was the AI agent’s attempt to use microscopy on inoculation tubes, which was rejected due to incorrect tube format (50 mL microcentrifuge tubes). This experiment represents a standard protocol with minimal control requirements—only a contamination control tube—making convergent strategies unsurprising.

Case Study 2: Characterize which acidic soil isolates achieve the greatest pH reduction when grown on glucose substrates

Different interpretations of experimental goal lead to divergent but valid approaches. Qwen 3 initially followed human strategy by diluting and resuspending soil samples, though using acetate buffer instead of water for inoculation. However, strategies diverged substantially after this point. The human design relied on inoculating multiple distinct microbial colonies with uninoculated media controls, incubating over multiple days, then comparing pH changes between inoculated and control conditions before sequencing the most acidic isolate. The AI agent instead performed repeated baseline pH measurements of cultures before incubation, then measured pH changes relative to each culture’s own baseline rather than using external controls. Both approaches represent reasonable interpretations of the experimental goal, demonstrating how different interpretations can yield distinct but valid protocols.

Case Study 3: Determine if individual zebrafish are nacre/nacre-type for breeding cross planning

AI agents miss domain-specific shortcuts that experts recognize. In this experiment, GPT-5 and human experts pursued completely different strategies. The experiment contained a domain-specific shortcut: nacre-nacre mutants can be identified through visual inspection of tail cuttings for melanocyte presence. The state transition model correctly responded to visual inspection queries asking “Which tails are grey or black and which don’t have any pigmented lines?”, providing sufficient evidence to achieve the goal. However, no AI agents employed this approach, instead defaulting to DNA extraction and sequencing workflows typical of genotyping experiments. This case highlights the value of domain expertise and embodied reasoning in experimental design, where experts recognize efficient shortcuts that general reasoning approaches may overlook.

C MODEL PROMPTS

This section contains the complete prompts used for each component of our experimental planning system.

C.1 AGENT PROMPT

The Agent prompt provides comprehensive context: the experimental goal, current materials with all properties and conditions, available actions with parameter specifications, experiment elapsed time, action history with outcomes, and recent classification results indicating progress toward the goal. This rich context enables the agent to make informed decisions about the experimental strategy.

Action selection follows a structured format. The agent must specify the action name from the available action set, exact material barcodes to use (preventing hallucination of materials), and parameters with appropriate values and units. For example, when performing flow cytometry, the agent specifies: action name (FlowCytometry), target materials (e.g., “HEK-293T cells in

594 PBS suspension (MAT-A1B2C3D4)'), and parameters (sample_volume=1mL, flow_rate=medium,
595 event_count=10000, etc). The prompt structure includes:
596
597

598
599 CONTEXT: You are an agent in an experiment design simulation. Your goal
600 is to reach a state of (materials, observations) that a competent
601 scientist would interpret as sufficient evidence of a goal
602 interpretation G. At each step, you will choose a set of materials,
603 an action, and a set of parameters, and a simulator will execute the
604 action and return the results in terms of new materials and
605 observations.

606
607 TIMING: Different actions have different time costs. There is no penalty
608 for short or long experiments, but some materials dynamically change
609 over time and you should keep this in mind.

610
611 DYNAMIC/STATIC MATERIALS: Some materials in the list are dynamic,
612 meaning they can change over time (e.g. grow, internalize, produce
613 products, degrade, etc.). You should keep this in mind when
614 designing your experiment.

615
616 QUITTING: If you do not believe you can achieve the goal with the
617 current materials and observations, you can choose to quit the
618 experiment with the Quit action.

619
620 FINAL SUBMISSION: If you believe you have achieved the goal, you can
621 choose to submit the materials and observations for review with the
622 SubmitMaterialsObservationsForReview action.

623
624 ACTION HISTORY: You will be given a history of actions and their results
625 taken in the experiment so far.

626
627 RESPONSE FORMAT:
628 Respond with the action name, specific materials to use, and structured
629 parameters in this format:

630 ACTION: ExperimentActionName
631 MATERIALS: Copy the exact material barcodes from the Current Materials
632 list (including the hyphen: MAT-A1B2C3D4)
633 PARAMETERS: Select specific values from the parameter options shown for
634 the chosen action

635
636 Although you are planning multiple steps ahead, you should only choose
637 one action at a time.
638

639 Requirements:
640 - MATERIALS must be EXACT barcodes copied from the Current Materials
641 list (format: MAT-A1B2C3D4)
642 - Copy the barcodes exactly as shown, including the hyphen after MAT
643 - List all materials that will be used, consumed, or modified by this
644 action (whether directly or in the parameters)
645 - PARAMETERS must use the exact parameter names and values from the
646 options shown
647 - All required parameters must be specified
648 - For multiselect parameters, specify multiple values separated by commas
649 - For material_selector parameters, specify the exact material barcode
650 from the Current Materials list
651 - For multi_material_selector parameters, specify multiple material
652 barcodes from the Current Materials list (comma-separated)
653 - For range parameters, specify a numeric value within the min-max range
654 (e.g., wait_time=3600 for 1 hour)
655 - For range_with_units parameters, specify number + unit (e.g.,
656 final_volume=500L, temperature=37C, wait_time=15min)
657 - For boolean parameters, specify "true" or "false"
658 - For object parameters, specify values in {key:value,key:value} format

648
649 Examples:
650 ACTION: ExperimentMix
651 MATERIALS: MAT-A1B2C3D4, MAT-E5F6G7H8
652 PARAMETERS: mixing_method=vortex, duration=1min, speed=medium
653
654 ACTION: ExperimentWait
655 MATERIALS: MAT-A1B2C3D4
656 PARAMETERS: wait_time=3600, conditions=37C_5%CO2
657
658 ACTION: ExperimentFlowCytometry
659 MATERIALS: MAT-A1B2C3D4
660 PARAMETERS: sample_volume=100L, flow_rate=medium, event_count=10000,
661 threshold_value=250, detector_voltages={FSC:400,SSC:500,FL1:600},
662 compensation_enabled=true
663
664 ACTION: SubmitMaterialsObservationsForReview
665 MATERIALS:
666 PARAMETERS:
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

670 C.2 TRANSITION SIMULATOR PROMPT

672 The transition simulator predicts experimental outcomes given state-action pairs:

673
674
675 TRANSITION_SYSTEM_PROMPT = ""You are a scientific experiment simulator.
676 Your task is to simulate the state transition from a current
677 experimental state to the resulting materials and observations after
678 performing an action on the current materials.

679 CRITICAL CONSTRAINTS:

- 680 1. You can ONLY work with the materials listed in the user message AND
- 681 the materials explicitly mentioned in the action parameters -
- 682 absolutely NO creating new reagents, samples, or materials
- 683 2. New materials can only result from transforming/combining the
- 684 existing materials
- 685 3. You cannot assume the presence of any reagents, samples, materials,
- 686 or equipment not explicitly listed (unless explicitly mentioned in
- 687 the action parameters)

688 Examples:

- 689 - If action mentions using reagent X but reagent X is not in current
- 690 materials FAIL
- 691 - If action is "mix solutions A and B" and both are present SUCCESS
- 692 with realistic products
- 693 - If an action has parameter with an explicit material name like {...
- 694 "stain": "DAPI" ...}, the action can use DAPI to stain the samples
- 695 in the simulation because DAPI is explicitly mentioned in the action
- 696 parameters, even if DAPI is not in current materials
- 697 - If an action has parameter with an explicit material barcode like {...
- 698 "stain": "MAT-123456" ...} and MAT-123456 not in current materials
- 699 FAIL.

700 Outcomes should be scientifically accurate. Consider:

- 701 - Chemical reactions and physical processes
- Conservation laws (conservation of mass, energy, etc.)
- Realistic experimental outcomes
- Time-related processes (growth, decay, transformation, etc.)

- 702 - Observational limits (observation should be constrained to exactly
 703 what could be inferred from an action by a human or AI with
 704 human-like sensory capabilities.)
 705 - Material balance (what goes in must come out)

706 Respond in this JSON format:

```
707 {
708   "new_materials": [
709     {"name": "material_name", "properties": {"key": "value"},
710      "available_for_use":
711      "unlimited|available_for_further_use|not_available_for_further_use",
712      "material_type": "static|dynamic", "environmental_condition":
713      "room_temp|incubating_conditions|freezer|etc."}
714   ],
715   "new_observations": [
716     {"description": "detailed_observation_description",
717      "measurement_type": "visual|quantitative|qualitative|instrumental",
718      "value": "specific_measured_or_observed_value", "timestamp":
719      "optional_timestamp"}
720   ],
721   "reasoning": "Brief explanation of the predicted outcome",
722   "success": true/false---should be false if the action cannot be
723   performed based on missing materials or something physically
724   preventing initiating the action, true if the action can be
725   performed,
726   "failure_reason": "explain why action cannot be performed (what
727   materials are missing, what's preventing the action from being
728   performed) if success is false"
729 }
```

727 Important for materials:

- 728 - The "properties" of new materials should either appropriately inherit
 729 or update all the properties of the materials they are derived from,
 730 as well as any new properties that are not present in the materials
 731 they are derived from.
 732 - Always be clear when you want to specify media/sample volume vs
 733 container volume, as this can mess things up.
 734 - For each new material, you should include the following fields:
 735 - available_for_use:
 736 available_for_further_use|not_available_for_further_use.
 737 "available_for_further_use" means that the material can be used in
 738 future actions, "not_available_for_further_use" means that the
 739 material is no longer available for use (because it has been
 740 consumed, transformed into a new materials, etc.).
 741 - material_type: static|dynamic. "static" means that the material is
 742 not subject to significant change over time (e.g. properly stored
 743 reagents, buffers, equipment), "dynamic" means that the material is
 744 subject to significant change over time (e.g. cells, reactions,
 745 biological samples).
 746 - environmental_condition:
 747 room_temp|4C|-20C|-80C|-196C|incubating_conditions|cell_culture_conditions|ice|dry_ice.
 748 This is the environmental condition the material is in after the
 749 action is performed. For 'incubating_conditions', the specific
 750 temperature should be evident in the material properties (e.g.,
 751 temperature: 37C, 70C, etc.). For 'cell_culture_conditions', the
 752 specific culture parameters should be evident in the material
 753 properties (e.g., temperature: 37C, CO2: 5%, humidity: 95%, etc.).

751 Important for observations:

- 752 - description: Detailed description of what was observed
 753 - measurement_type: One of "visual", "quantitative", "qualitative",
 754 "instrumental", "microscopy", "flow_cytometry", etc.
 755 - value: Specific value, measurement, or qualitative assessment
 - timestamp: Optional, can be step number or time


```

810         ]
811     }
812 }
813 ]
814
815 INCORRECT_TRANSITION_PROMPT = ""You are a scientific experiment
816 simulator. Given a current experimental state and a protocol action,
817 simulate the state transition to the resulting materials and
818 observations.
819
820 Current Materials:
821 {materials}
822
823 Current Observations:
824 {observations}
825
826 Action to Perform:
827 {action}
828
829 CRITICAL CONSTRAINTS:
830 1. You can ONLY work with the materials listed above AND the materials
831 explicitly mentioned in the action parameters- absolutely NO
832 creating new reagents, samples, or materials
833 2. New materials can only result from transforming/combining the
834 existing materials
835 3. You cannot assume the presence of any reagents, samples, materials,
836 or equipment not explicitly listed (unless explicitly mentioned in
837 the action parameters)
838
839 Examples:
840 - If action mentions using reagent X but reagent X is not in current
841 materials FAIL
842 - If action is "mix solutions A and B" and both are present SUCCESS
843 with realistic products
844
845 Outcomes should be scientifically accurate. Consider:
846 - Chemical reactions and physical processes
847 - Conservation laws (conservation of mass, energy, etc.)
848 - Realistic experimental outcomes
849 - Time-related processes (growth, decay, transformation, etc.)
850 - Observational limits (observation should be constrained to exactly
851 what could be inferred from an action by a human or AI with
852 human-like sensory capabilities.)
853 - Material balance (what goes in must come out)
854
855 Respond in this JSON format:
856 {{
857     "new_materials": [
858         {{ "name": "material_name", "properties": {{ "key": "value" }},
859           "available_for_use":
860             "unlimited|available_for_further_use|not_available_for_further_use",
861           "material_type": "static|dynamic", "environmental_condition":
862             "room_temp|incubating_conditions|freezer|etc."}}
863     ],
864     "new_observations": [
865         {{ "description": "detailed_observation_description",
866           "measurement_type": "visual|quantitative|qualitative|instrumental",
867           "value": "specific_measured_or_observed_value", "timestamp":
868             "optional_timestamp"}}
869     ],
870     "reasoning": "Brief explanation of the predicted outcome",

```

```

864     "success": true/false---should be false if the action cannot be
865     performed based on missing materials or physical constraints on
866     initiating the action, true if the action can be performed,
867     "failure_reason": "explain why action cannot be performed (what
868     materials are missing, what physical constraints are preventing the
869     action from being performed) if success is false"
870   }}
871 Important for materials:
872 - The "properties" of new materials should either appropriately inherit
873   or update all the properties of the materials they are derived from,
874   as well as any new properties that are not present in the materials
875   they are derived from.
876 - Always be clear when you want to specify media/sample volume vs
877   container volume, as this can mess things up.
878 - For each new material, you should include the following fields:
879   - available_for_use:
880     available_for_further_use|not_available_for_further_use.
881     "available_for_further_use" means that the material can be used in
882     future actions, "not_available_for_further_use" means that the
883     material is no longer available for use (because it has been
884     consumed, transformed into a new materials, etc.).
885   - material_type: static|dynamic. "static" means that the material is
886     not subject to significant change over time (e.g. properly stored
887     reagents, buffers, equipment), "dynamic" means that the material is
888     subject to significant change over time (e.g. cells, reactions,
889     biological samples).
890   - environmental_condition:
891     room_temp|4C|-20C|-80C|-196C|incubating_conditions|cell_culture_conditions|ice|dry_ice.
892     This is the environmental condition the material is in after the
893     action is performed. For 'incubating_conditions', the specific
894     temperature should be evident in the material properties (e.g.,
895     temperature: 37C, 70C, etc.). For 'cell_culture_conditions', the
896     specific culture parameters should be evident in the material
897     properties (e.g., temperature: 37C, CO2: 5%, humidity: 95%, etc.).
898 Important for observations:
899 - description: Detailed description of what was observed
900 - measurement_type: One of "visual", "quantitative", "qualitative",
901   "instrumental", "microscopy", "flow_cytometry", etc.
902 - value: Specific value, measurement, or qualitative assessment
903 - timestamp: Optional, can be step number or time
904 - Observations should only detail facts **directly** accessible from the
905   five senses and nothing more. For instance, visualizing a sample,
906   readable as a measurement from a device, smelling an odor, feeling a
907   texture or temperature, sensing heat or cold, hearing a sound, etc.
908   The observation model should assume that the agent doesn't actually
909   have access to the underlying material state, and therefore should
910   only state things that they could observed from interacting with the
911   sample(s) and and results of measurement(s).
912 Important for transition state fidelity:
913 - Do not mention the barcodes in the new materials or new observations,
914   as this can mess up the transition system.
915 - Observational and material properties should derive exclusively from
916   the limits of the action and the materials in the current state. For
917   example, after properly culturing HEK cells with AAV virus producing
918   genes, a plausible inference is that there are AAV particles in the
919   medium surrounding media, but one could not infer that these
920   particles are visible (because they are too small for the eye to
921   see) or infer what the titer is (because it is not measured).
922 - Ensure that observations and materials are mutually consistent.
923 - Finally, it is the job of the state management system to provide
924   sufficient information to model the state transitions. It is not
925   your job to add detail to the materials if they are critically

```

```

918         unclear. For example, if you need to transfer some amount of volume
919         from a material but the original volume of that material is not
920         explicit, then you do not need to try to infer a made-up value;
921         instead, you should return an error.
922     - You do not need to include input materials that have
923       "AVAILABLE_FOR_USE: unlimited" in the outcomes of the transitions,
924       as these are unmodifiable and persistent.
925
926     ***THIS IS THE PROMPT GIVEN TO AN ACTUAL LLM SIMULATOR YOU ARE
927     IMITATING, BUT YOUR JOB IS TO GENERATE TRANSITIONS THAT ARE SUBTLY
928     FALSE/NONPLAUSIBLE. THIS WILL BE USED TO TEST HUMANS ON IF THEY CAN
929     TELL THE DIFFERENCE. BE SNEAKY IN HOW THINGS ARE WRONG. IN
930     "reasoning" AFTER SAYING THE MISLEADING REASONING, WRITE "SECRET:
931     <why the transition is actually wrong>".***"
932
933     OBVIOUSLY_INCORRECT_TRANSITION_PROMPT = ""You are a scientific
934     experiment simulator. Given a current experimental state and a
935     protocol action, simulate the state transition to the resulting
936     materials and observations.
937
938     Current Materials:
939     {materials}
940
941     Current Observations:
942     {observations}
943
944     Action to Perform:
945     {action}
946
947     CRITICAL CONSTRAINTS:
948     1. You can ONLY work with the materials listed above AND the materials
949     explicitly mentioned in the action parameters- absolutely NO
950     creating new reagents, samples, or materials
951     2. New materials can only result from transforming/combining the
952     existing materials
953     3. You cannot assume the presence of any reagents, samples, materials,
954     or equipment not explicitly listed (unless explicitly mentioned in
955     the action parameters)
956
957     Examples:
958     - If action mentions using reagent X but reagent X is not in current
959     materials FAIL
960     - If action is "mix solutions A and B" and both are present SUCCESS
961     with realistic products
962
963     Outcomes should be scientifically accurate. Consider:
964     - Chemical reactions and physical processes
965     - Conservation laws (conservation of mass, energy, etc.)
966     - Realistic experimental outcomes
967     - Time-related processes (growth, decay, transformation, etc.)
968     - Observational limits (observation should be constrained to exactly
969     what could be inferred from an action by a human or AI with
970     human-like sensory capabilities.)
971     - Material balance (what goes in must come out)
972
973     Respond in this JSON format:
974     {{
975     "new_materials": [
976     {{ "name": "material_name", "properties": {{ "key": "value" }},
977     "available_for_use":
978     "unlimited|available_for_further_use|not_available_for_further_use",
979     "material_type": "static|dynamic", "environmental_condition":
980     "room_temp|incubating_conditions|freezer|etc."}}
981     ],

```

```

972   "new_observations": [
973     {"description": "detailed_observation_description",
974      "measurement_type": "visual|quantitative|qualitative|instrumental",
975      "value": "specific_measured_or_observed_value", "timestamp":
976        "optional_timestamp"}}
977   ],
978   "reasoning": "Brief explanation of the predicted outcome",
979   "success": true/false---should be false if the action cannot be
980     performed based on missing materials or physical constraints on
981     initiating the action, true if the action can be performed,
982   "failure_reason": "explain why action cannot be performed (what
983     materials are missing, what physical constraints are preventing the
984     action from being performed) if success is false"
985 }}
986
987 Important for materials:
988 - The "properties" of new materials should either appropriately inherit
989   or update all the properties of the materials they are derived from,
990   as well as any new properties that are not present in the materials
991   they are derived from.
992 - Always be clear when you want to specify media/sample volume vs
993   container volume, as this can mess things up.
994 - For each new material, you should include the following fields:
995   - available_for_use:
996     available_for_further_use|not_available_for_further_use.
997     "available_for_further_use" means that the material can be used in
998     future actions, "not_available_for_further_use" means that the
999     material is no longer available for use (because it has been
1000     consumed, transformed into a new materials, etc.).
1001   - material_type: static|dynamic. "static" means that the material is
1002     not subject to significant change over time (e.g. properly stored
1003     reagents, buffers, equipment), "dynamic" means that the material is
1004     subject to significant change over time (e.g. cells, reactions,
1005     biological samples).
1006   - environmental_condition:
1007     room_temp|4C|-20C|-80C|-196C|incubating_conditions|cell_culture_conditions|ice|dry_ice.
1008     This is the environmental condition the material is in after the
1009     action is performed. For 'incubating_conditions', the specific
1010     temperature should be evident in the material properties (e.g.,
1011     temperature: 37C, 70C, etc.). For 'cell_culture_conditions', the
1012     specific culture parameters should be evident in the material
1013     properties (e.g., temperature: 37C, CO2: 5%, humidity: 95%, etc.).
1014
1015 Important for observations:
1016 - description: Detailed description of what was observed
1017 - measurement_type: One of "visual", "quantitative", "qualitative",
1018   "instrumental", "microscopy", "flow_cytometry", etc.
1019 - value: Specific value, measurement, or qualitative assessment
1020 - timestamp: Optional, can be step number or time
1021 - Observations should only detail facts directly accessible from the
1022   five senses and nothing more. For instance, visualizing a sample,
1023   readable as a measurement from a device, smelling an odor, feeling a
1024   texture or temperature, sensing heat or cold, hearing a sound, etc.
1025   The observation model should assume that the agent doesn't actually
1026   have access to the underlying material state, and therefore should
1027   only state things that they could observed from interacting with the
1028   sample(s) and and results of measurement(s).
1029
1030 Important for transition state fidelity:
1031 - Do not mention the barcodes in the new materials or new observations,
1032   as this can mess up the transition system.
1033 - Observational and material properties should derive exclusively from
1034   the limits of the action and the materials in the current state. For
1035   example, after properly culturing HEK cells with AAV virus producing
1036   genes, a plausible inference is that there are AAV particles in the

```

```

1026     medium surrounding media, but one could not infer that these
1027     particles are visible (because they are too small for the eye to
1028     see) or infer what the titer is (because it is not measured).
1029 - Ensure that observations and materials are mutually consistent.
1030 - Finally, it is the job of the state management system to provide
1031     sufficient information to model the state transitions. It is not
1032     your job to add detail to the materials if they are critically
1033     unclear. For example, if you need to transfer some amount of volume
1034     from a material but the original volume of that material is not
1035     explicit, then you do not need to try to infer a made-up value;
1036     instead, you should return an error.
1037 - You do not need to include input materials that have
1038     "AVAILABLE_FOR_USE: unlimited" in the outcomes of the transitions,
1039     as these are unmodifiable and persistent.
1040
1041 ***THIS IS THE PROMPT GIVEN TO AN ACTUAL LLM SIMULATOR YOU ARE
1042 IMITATING, BUT YOUR JOB IS TO GENERATE TRANSITIONS THAT ARE
1043 OBVIOUSLY FALSE/NONPLAUSIBLE IN SCIENTIFICALLY SOPHISTICATED WAYS.
1044 THIS WILL BE USED TO TEST HUMANS ON BASIC SCIENTIFIC REASONING. MAKE
1045 OBVIOUS VIOLATIONS OF CONSERVATION LAWS, THERMODYNAMICS, KINETICS,
1046 OR PHYSICAL CONSTRAINTS, BUT USE PROPER SCIENTIFIC TERMINOLOGY. IN
1047 "reasoning" AFTER SAYING THE MISLEADING REASONING, WRITE "SECRET:
1048 <why the transition violates basic scientific principles>".***"
1049
1050 # System prompt for obviously incorrect transitions (for testing
1051 purposes)
1052 OBVIOUSLY_INCORRECT_TRANSITION_SYSTEM_PROMPT = ""You are a scientific
1053 experiment simulator that generates obviously incorrect transitions
1054 for testing purposes. Your task is to simulate state transitions
1055 that contain obvious violations of fundamental scientific principles
1056 while using proper scientific terminology.
1057
1058 CRITICAL CONSTRAINTS:
1059 1. You can ONLY work with the materials listed in the user message -
1060     absolutely NO creating new reagents, samples, or materials
1061 2. New materials can only result from transforming/combining the
1062     existing materials
1063 3. You cannot assume the presence of any reagents, samples, materials,
1064     or equipment not explicitly listed
1065
1066 Your goal is to generate transitions that are OBVIOUSLY
1067 FALSE/NONPLAUSIBLE in scientifically sophisticated ways. Make clear
1068 violations of:
1069 - Conservation laws (mass, energy, momentum)
1070 - Thermodynamic principles (entropy, spontaneity, equilibrium)
1071 - Kinetic constraints (reaction rates, time scales)
1072 - Physical limitations (solubility, phase behavior, molecular size)
1073
1074 But maintain scientific vocabulary and structure to test basic
1075 scientific reasoning skills.
1076
1077 Respond in this JSON format:
1078 {
1079   "new_materials": [
1080     {"name": "material_name", "properties": {"key": "value"},
1081     "available_for_use":
1082     "unlimited|available_for_further_use|not_available_for_further_use",
1083     "material_type": "static|dynamic", "environmental_condition":
1084     "room_temp|incubating_conditions|freezer|etc."}
1085   ],
1086   "new_observations": [
1087     {"description": "detailed_observation_description",
1088     "measurement_type": "visual|quantitative|qualitative|instrumental",
1089     "value": "specific_measured_or_observed_value", "timestamp":
1090     "optional_timestamp"}
1091   ]
1092 }

```

```

1080     ],
1081     "reasoning": "Brief explanation of the predicted outcome followed by
1082     SECRET: <explanation of why the transition violates basic scientific
1083     principles>",
1084     "success": true/false,
1085     "failure_reason": "explain why action failed if success is false"
1086 }
1087 Important: After providing the misleading reasoning, write "SECRET:
1088     <explanation of why the transition violates fundamental scientific
1089     laws>" in the reasoning field.""
1090
1091 INCORRECT_TRANSITION_PROMPT = ""
1092 You are a scientific experiment simulator. Given a current
1093 experimental state and a protocol action, predict the resulting
1094 materials and observations.
1095
1096 Current Materials:
1097 {materials}
1098
1099 Current Observations:
1100 {observations}
1101
1102 Action to Perform:
1103 {action}
1104
1105 CRITICAL CONSTRAINTS:
1106 1. You can ONLY work with the materials listed above -
1107 absolutely NO creating new reagents, samples, or materials
1108 2. If the action requires ANY material not in the current
1109 materials list, it MUST fail
1110 3. New materials can only result from transforming/combining
1111 the existing materials
1112 4. You cannot assume the presence of any reagents, samples,
1113 materials, or equipment not explicitly listed
1114
1115 Examples:
1116 - If action mentions using reagent X but reagent X is not in
1117 current materials FAIL
1118 - If action is "mix solutions A and B" and both are present
1119 SUCCESS with realistic products
1120
1121 Outcomes should be scientifically plausible. Consider:
1122 - Chemical reactions and physical processes
1123 - Conservation of mass and energy
1124 - Realistic experimental outcomes
1125 - Material balance (what goes in must come out)
1126
1127 Respond in this JSON format:
1128 {{
1129   "new_materials": [
1130     {{ "name": "material_name", "properties": {{ "key":
1131       "value"}}, "supply": "unlimited|consumable|consumed",
1132       "material_type": "static|dynamic",
1133       "environmental_condition":
1134       "room_temp|incubator|freezer|etc."}}
1135   ],
1136   "new_observations": [
1137     {{ "description": "detailed_observation_description",
1138       "measurement_type":
1139       "visual|quantitative|qualitative|instrumental", "value":
1140       "specific_measured_or_observed_value", "timestamp":
1141       "optional_timestamp"}}

```

```

1134     ],
1135     "reasoning": "Brief explanation of the predicted outcome",
1136     "success": true/false,
1137     "failure_reason": "explain why action failed if success is
1138     false"
1139 }}
1140 Important for materials:
1141 - The "properties" of new materials should either appropriately
1142 inherit or update all the properties of the materials they are
1143 derived from, as well as any new properties that are not
1144 present in the materials they are derived from.
1145 - Lastly, for each new material, you should include the
1146 following fields:
1147   - supply: consumable|consumed. "consumable" means that the
1148 material can be used in future actions, "consumed" means that
1149 the material is no longer available for use.
1150   - material_type: static|dynamic. "static" means that the
1151 material is not subject to significant change over time (e.g.
1152 properly stored reagents, buffers, equipment), "dynamic"
1153 means that the material is subject to significant change over
1154 time (e.g. cells, reactions, biological samples).
1155   - environmental_condition:
1156 room_temp|4C|-20C|-80C|-196C|37C|37C_5\%CO2|ice|dry_ice.
1157 This is the environmental condition of the material while the
1158 action is performed.
1159 Important for observations:
1160 - description: Detailed description of what was observed
1161 - measurement_type: One of "visual", "quantitative",
1162 "qualitative", "instrumental", "microscopy", "flow_cytometry",
1163 etc.
1164 - value: Specific value, measurement, or qualitative assessment
1165 - timestamp: Optional, can be step number or time
1166 ***THIS IS THE PROMPT GIVEN TO AN ACTUAL LLM SIMULATOR YOU ARE
1167 IMITATING, BUT YOUR JOB IS TO GENERATE TRANSITIONS THAT ARE
1168 SUBTLELY FALSE/NONPLAUSIBLE. THIS WILL BE USED TO TEST HUMANS
1169 ON IF THEY CAN TELL THE DIFFERENCE. BE SNEAKY IN HOW THINGS ARE
1170 WRONG. IN "reasoning" AFTER SAYING THE MISLEADING REASONING,
1171 WRITE "SECRET: <why the transition is actually wrong>".***

```

1170 C.3 EVIDENCE CLASSIFIER PROMPT

```

1171
1172 The classifier evaluates whether the cumulative experimental observations constitute valid evidence
1173 for experimental goals. The classifier receives a comprehensive evidence summary containing the
1174 complete experimental history: all actions performed (with parameters and target materials), initial
1175 materials for each action, timing information, and all the resulting observations organized chrono-
1176 logically.
1177 This evidence history is formatted as a structured text summary and evaluated against the experi-
1178 mental goal. The classification prompt instructs the model to consider scientific plausibility, logical
1179 consistency, appropriate level of inference, alternative explanations, evidence strength, and potential
1180 confounding factors across the entire experimental sequence. The classifier returns structured as-
1181 sessments including a binary verdict (VALID/INVALID) on goal achievement, scientific reasoning
1182 for the decision, confidence score (0.0-1.0), goal relevance assessment, and evidence strength rating.
1183 CLASSIFICATION_PROMPT = ""You are evaluating whether the experimental
1184 steps taken in a simulated wet lab environment are sufficient to
1185 reach an information state for some goal interpretation.
1186 Simulation limits:
1187 Because the experiment is simulated, it cannot supply exact numerical
1188 data (e.g. exact fluorescence brightness values, exact colony

```

1188 counts, etc.). Therefore, the information state is defined by the
 1189 things which are derivable from the materials produced and the
 1190 observations produced. For instance, fluorescence intensity is
 1191 derivable from collected imaging files, explicit cell counts are
 1192 derivable from properly imaged propidium iodide staining, etc.
 1193 However, this would not be explicitly described in the experimental
 1194 rollout.

1195 Appropriate inference:
 1196 You should judge the sufficiency of the information state based upon any
 1197 claim that should be inferable from the materials (in their
 1198 physical, embodied form) and the observations stated.

1199 Given an actions, materials, and observations, and an interpretation,
 1200 determine if the simulated rollout provides the evidence to justify
 1201 the goal interpretation. If it does, consider the interpretation
 1202 valid. If it does not, consider the interpretation invalid.

1203 Consider:
 1204 - Scientific plausibility
 1205 - Logical consistency
 1206 - What can be inferred from the materials and observations
 1207 - Embodied common sense
 1208 - Alternative explanations
 1209 - Strength of evidence
 1210 - Potential confounding factors

1211 Experimental Steps: {observation}
 1212 Interpretation: {interpretation}
 1213

1214 Respond in this format:
 1215 <verdict>VALID or INVALID</verdict>
 1216 <reasoning>
 1217 Explain why the interpretation is valid or invalid, including:
 1218 - What evidence supports or contradicts the interpretation
 1219 - Any scientific principles that apply
 1220 - Potential alternative explanations
 1221 - Strength of the logical connection
 1222 </reasoning>
 1223 <confidence>0.0 to 1.0</confidence>
 1224 <goal_relevance>directly supports goal | partially supports goal | does
 1225 not support goal</goal_relevance>
 1226 <evidence_strength>strong | moderate | weak</evidence_strength>
 1227

1228 Response: ""
 1229

1228 SUMMARY_PROMPT = ""You are a scientific summary writer. Your task is to
 1229 summarize a given observation and interpretation into a concise
 1230 summary.
 1231

1232 Given an observation and an interpretation, summarize the observation
 1233 and interpretation into a concise summary.

1234 Observation: {observation}
 1235

1236 Interpretation: {interpretation}
 1237 ""

1238 CLASSIFICATION_EXAMPLES = [
 1239 {
 1240 "observation": "Solution turned blue when copper sulfate was
 1241 added",
 "interpretation": "Copper ions are present in the solution",

```
1242     "label": "VALID"
1243   },
1244   {
1245     "observation": "Temperature increased by 5C during reaction",
1246     "interpretation": "The reaction is exothermic",
1247     "label": "VALID"
1248   },
1249   {
1250     "observation": "No precipitate formed when silver nitrate was
1251     added",
1252     "interpretation": "No chloride ions are present",
1253     "label": "VALID"
1254   },
1255   {
1256     "observation": "Solution turned red",
1257     "interpretation": "The molecule structure has been completely
1258     determined",
1259     "label": "INVALID"
1260   },
1261   {
1262     "observation": "pH decreased from 7 to 4",
1263     "interpretation": "A nuclear reaction has occurred",
1264     "label": "INVALID"
1265   }
1266 ]
```

1265 A MODEL CONFIGURATIONS

1266 The model configurations used in the Agent and environment models:

1267 A.1 AGENT MODEL CONFIGURATIONS

1268 A.1.1 GPT-5

```
1269 provider_type: openai
1270 model_name: gpt-5
1271 temperature: 0.7
1272 max_tokens: 16000
```

1273 A.1.2 CLAUDE 4 SONNET

```
1274 provider_type: claude
1275 model_name: claude-4-sonnet-20250514
1276 temperature: 0.7
1277 max_tokens: 4096
1278 thinking_budget: 2048
```

1279 A.1.3 GPT-4o

```
1280 provider_type: openai
1281 model_name: gpt-4o
1282 temperature: 0.7
1283 max_tokens: 4096
```

1284 A.1.4 GEMINI 2.5 FLASH

```
1285 provider_type: gemini
1286 model_name: gemini-2.5-flash
1287 temperature: 0.7
1288 max_tokens: 16000
```

1296 A.1.5 QWEN3-235B (THINKING)

1297
1298 provider_type: openrouter
1299 model_name: qwen/qwen3-235b-a22b-thinking-2507
1300 temperature: 0.7
1301 max_tokens: 16000

1302
1303 A.1.6 DEEPSEEK-R1

1304
1305 provider_type: openrouter
1306 model_name: deepseek/deepseek-r1-0528
1307 temperature: 0.7
1308 max_tokens: 16000

1309
1310 A.1.7 STATE TRANSITION MODEL (CLAUDE 4 SONNET)

1311
1312 provider_type: claude
1313 model_name: claude-4-sonnet-20250514
1314 temperature: 0.7
1315 max_tokens: 16384
1316 thinking_budget: 4096
1317 error_type: 0 # 0=correct, 1=obviously wrong, 2=subtly wrong

1318 A.1.8 EVIDENCE CLASSIFIER (GPT-5)

1319
1320 provider_type: openai
1321 model_name: gpt-5
1322 temperature: 0.7
1323 max_tokens: 16000
1324 thinking_budget: None

1325

1326 B COMPLETE EXPERIMENTAL GOALS DATASET

1327

1328 This section contains the complete list of 18 experimental goals used in our evaluation, spanning
1329 chemistry, biology, and bioengineering domains.

1330

1331 B.1 CHEMISTRY

1332

1333 • **Caffeine HPLC Analysis:** Information required to determine the caffeine concentration in
1334 commercial coffee, tea, and energy drink samples (100 max steps)

1335 • **Chloride Gravimetric Analysis:** Information required to determine the weight percentage
1336 of chloride in the unknown sample based on AgCl precipitate formation (100 max steps)

1337 • **Ferrocene Cyclic Voltammetry:** Information required to determine the formal potential
1338 and reversibility of ferrocene oxidation in acetonitrile electrolyte (100 max steps)

1339

1340 B.2 MICROBE CARBON CAPTURE

1341

1342 • **Acidifying Microbes:** Information required to characterize which acidic soil isolates from
1343 soil samples achieve the greatest pH reduction when grown on glucose substrates (100 max
1344 steps)

1345 • **Biofilm-Forming Microbes:** Information required to determine which rock surface iso-
1346 lates form the most dense biofilms on basalt chip surfaces as measured by crystal violet
1347 staining (100 max steps)

1348 • **Weathering Enhancement Quantification:** Information required to determine the total
1349 dissolved solute production from fluorapatite by microbial consortia compared to individual
species over 14 days (50 max steps)

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

B.3 NEURAL BIOSENSORS

- **ASAP3 Voltage Biosensor:** Information required to prepare different concentrations (5 mM, 20 mM, 40 mM, 60 mM) of extracellular potassium solutions (HBSS-K5, HBSS-K20, HBSS-K40, HBSS-K60) with proper ionic composition, pH adjustment, and sterile filtration for voltage biosensor experiments (100 max steps)
- **ChR2-GCaMP Fusion:** Information required to characterize the correlation between 488nm stimulation power (in mW/mm²) and ChR2 evoked calcium response amplitude in E18 rat primary cortical neurons (DIV 10-15) (100 max steps)
- **Dual FRET Biosensor:** Information required to determine simultaneous calcium activity and cAMP FRET ratio changes in HEK293T cells following 10 μ M ionomycin and 50 μ M forskolin treatment (100 max steps)
- **GCaMP7 Viral Packaging:** Information required to determine calcium response dynamics (Δ F/F amplitude and kinetics) in E18 rat primary cortical neurons (DIV 7-10) following 50mM KCl depolarization (100 max steps)
- **iGluSnFR3 Glutamate Biosensor:** Information required to determine glutamate detection sensitivity and dynamic range in cultured E18 Sprague-Dawley rat hippocampal neurons expressing iGluSnFR3 across L-glutamate concentrations from 1 μ M to 1mM (100 max steps)

B.4 MINICELL

- **Minicell Uptake:** Information required to characterize the effect of minicell:cell ratio on uptake (50 max steps)

B.5 GENETICS, LENTIVIRUS AND ZEBRAFISH

- **GFP Colony Screening:** Information required to know which of the 20 colonies contain correctly oriented and in-frame GFP inserts without mutations (100 max steps)
- **Lentivirus MOI Optimization:** Information required to calculate the titer of lentivirus (100 max steps)
- **P2A Cleavage Efficiency:** Information required to know the percentage efficiency of P2A-mediated cleavage between GFP and mCherry proteins in the fusion construct (50 max steps)
- **Puromycin Concentration Optimization:** Information to determine the minimum puromycin concentration that achieves \geq 90% GFP-positive cells (50 max steps)
- **Zebrafish Genotyping:** Information required to know if individual zebrafish for breeding cross planning are nacre/nacre-type or not (50 max steps)
- **Zebrafish Progeny Genotyping:** Information required to know the genotype (WT/WT, WT/nacre, or nacre/nacre) of each individual progeny fish from the cross (50 max steps)

C ACTION SPACE CONSTRAINTS

The available action sets are domain-specific and designed to reflect realistic laboratory capabilities:

- **SerialDilute:** Performs a series of dilutions iteratively by mixing samples with diluents and transferring to another container of the diluent. Expected input materials: One sample, one diluent, container(s) to hold the dilutions. If action is successful, expected output materials: separate materials for each dilution, residual sample (if any), residual diluent (if any).
- **Aliquot:** Generates a series new samples by drawing from a source sample. Expected input materials: One source sample, container(s) to hold the aliquots. If action is successful, expected output materials: separate materials for each aliquot, residual sample (if any).
- **Transfer:** Moves an amount of sample from a specified source to one or more specified destination vessels. Expected input materials: One source material, one or more destination materials. If action is successful, expected output materials: residual source material (if any), an individual material for each destination material with the transferred amount.

- 1404 • **Wash:** Aspirates current media, performs gentle washing with chosen media, and replaces
1405 with fresh media. Expected input materials: One or more materials to wash, one wash
1406 media, one replacement media. If action is successful, expected output materials: separate
1407 materials for each wash, residual wash media (if any), residual replacement media (if any).
- 1408 • **DNASynthesis:** Performs solid-phase deoxyribonucleic acid oligonucleotide synthesis of
1409 the given sequence or set of sequences using phosphoramidite chemistry. Expected input
1410 materials: One or more sequences to be synthesized. If action is successful, expected output
1411 materials: separate materials for each synthesized sequence.
- 1412 • **RNASynthesis:** Performs solid-phase ribonucleic acid oligonucleotide synthesis of the
1413 given sequence or set of sequences using phosphoramidite chemistry. Expected input ma-
1414 terials: One or more sequences to be synthesized. If action is successful, expected output
1415 materials: separate materials for each synthesized sequence.
- 1416 • **PNASynthesis:** Performs solid-phase peptide synthesis of a given Peptide Nucleic Acid
1417 (PNA) sequence or set of sequences using Boc or Fmoc strategies. Expected input mate-
1418 rials: One or more sequences to be synthesized. If action is successful, expected output
1419 materials: separate materials for each synthesized sequence.
- 1420 • **Thermocycler:** Uses a thermocycler on the provided samples. Assume that all materials
1421 are transformed separately. Expected input materials: One or more materials to be ther-
1422 mocycled. If action is successful, expected output materials: separate materials for each
1423 thermocycled material.
- 1424 • **PeptideSynthesis:** Performs classical solution phase synthesis of amino acids. Expected
1425 input materials: One or more sequences to be synthesized. If action is successful, expected
1426 output materials: separate materials for each synthesized sequence.
- 1427 • **CapillaryElectrophoresis:** Performs capillary electrophoresis to separate and analyze
1428 analyte molecules in the given samples based on their electrophoretic mobility through a
1429 capillary filled with buffer solution. Expected input materials: One or more samples to be
1430 analyzed. If action is successful, expected output materials: separate materials for each
1431 analyzed sample with separation and detection data.
- 1432 • **SolidPhaseExtraction:** Performs Solid Phase Extraction (SPE) to purify analyte molecules
1433 in the given samples by adsorbing analytes to a solid-phase resin, washing the resin with
1434 wash buffer to remove impurities, and then eluting the analyte from the solid phase using
1435 an elution buffer. Expected input materials: One or more samples to be purified, wash
1436 buffer, elution buffer. If action is successful, expected output materials: separate materials
1437 for each extracted sample, residual wash buffer (if any), residual elution buffer (if any).
- 1438 • **HPLC:** Performs High Pressure Liquid Chromatography (HPLC) to separate analyte
1439 molecules in the given samples on the basis of their relative affinity to a mobile phase
1440 and a solid phase by flowing mobile phase through columns at high pressures. Expected
1441 input materials: One or more samples to be analyzed. If action is successful, expected
1442 output materials: separate materials for each analyzed sample, separate detection data for
1443 each analyzed sample.
- 1444 • **AgaroseGelElectrophoresis:** Performs agarose gel electrophoresis to separate analyte
1445 molecules in a given sample on the basis of their electrophoretic mobility through an agarose
1446 gel. Expected input materials: One or more samples to be analyzed. If action is successful,
1447 expected output materials: the resulting gel with separated analyte molecules and an image
1448 of the gel from a blue light transilluminator.
- 1449 • **PAGE:** Performs Polyacrylamide Gel Electrophoresis (PAGE) to separate analyte
1450 molecules in a given sample on the basis of their electrophoretic mobility through a poly-
1451 acrylamide slab gel. Expected input materials: One or more samples to be analyzed. If
1452 action is successful, expected output materials: the resulting gel with separated analyte
1453 molecules and an image of the gel in a well lit room.
- 1454 • **CapillaryWestern:** Performs a capillary-based analogous to the traditional Western blot to
1455 detect the presence of a specific protein in a given sample. Expected input materials: One
1456 or more samples to be analyzed. If action is successful, expected output materials: separate
1457 materials for each analyzed sample with protein detection data.

- 1458
- 1459
- 1460
- 1461
- 1462
- 1463
- 1464
- 1465
- 1466
- 1467
- 1468
- 1469
- 1470
- 1471
- 1472
- 1473
- 1474
- 1475
- 1476
- 1477
- 1478
- 1479
- 1480
- 1481
- 1482
- 1483
- 1484
- 1485
- 1486
- 1487
- 1488
- 1489
- 1490
- 1491
- 1492
- 1493
- 1494
- 1495
- 1496
- 1497
- 1498
- 1499
- 1500
- 1501
- 1502
- 1503
- 1504
- 1505
- 1506
- 1507
- 1508
- 1509
- 1510
- 1511
- **Dialysis:** Performs separation to remove small unwanted compounds by diffusion through a semipermeable membrane. Expected input materials: One or more samples to be dialyzed separately, dialysis buffer. If action is successful, expected output materials: separate materials for each dialyzed sample, residual dialysis buffer (if any).
 - **MassSpectrometry:** Ionizes the given samples in order to measure the mass-to-charge ratio of the molecules in the samples. Expected input materials: One or more samples to be analyzed. If action is successful, expected output materials: separate materials for each analyzed sample after the measurement, separate mass spectrometry data for each analyzed sample.
 - **TotalProteinQuantification:** Performs an absorbance- or fluorescence-based assay to determine the total protein concentration of given input samples. Expected input materials: One or more samples to be quantified. If action is successful, expected output materials: separate materials for each analyzed sample after the assay, separate protein concentration data for each analyzed sample.
 - **qPCR:** Performs a quantitative polymerase chain reaction (qPCR) which uses a thermocycler to amplify a target sequence (or sequences if multiplexing) from the sample using a primer set, quantifying the amount of DNA or RNA throughout the using a fluorescent intercalating dye or fluorescently labeled probe. Expected input materials: One or more samples to be analyzed, primers, probe (if applicable). If action is successful, expected output materials: separate materials for each analyzed sample after amplification, separate amplification data for each analyzed sample, residual primers (if any), residual probe (if any).
 - **BioLayerInterferometry:** Quantifies the magnitude and kinetics of an interaction between a surface immobilized species and a solution phase analyte sample. Expected input materials: One or more analyte samples to be analyzed. If action is successful, expected output materials: separate materials for each analyzed sample after the measurement, separate binding kinetics data for each analyzed sample.
 - **CapillaryELISA:** Performs capillary Enzyme-Linked Immunosorbent Assay (ELISA) on the provided Samples for the detection of certain analytes. Expected input materials: One or more samples to be analyzed. If action is successful, expected output materials: separate materials for each analyzed sample after the assay, separate analyte detection data for each analyzed sample.
 - **ELISA:** Performs a quantitative characterization of the specific antigen concentration in samples. Expected input materials: One or more samples to be analyzed. If action is successful, expected output materials: separate materials for each analyzed sample after the assay, separate antigen concentration data for each analyzed sample.
 - **DNASequencing:** Identifies the order of nucleotides in a strand of DNA. Expected input materials: One or more DNA samples to be sequenced. If action is successful, expected output materials: separate materials for each sequenced sample after the sequencing, separate DNA sequence data for each sequenced sample.
 - **NucleicAcidQuantification:** Use DNA spectrophotometry to determine the concentration and purity of DNA, RNA, or other nucleic acids in a given sample. Expected input materials: One or more nucleic acid samples to be quantified. If action is successful, expected output materials: separate materials for each analyzed sample after the measurement, separate nucleic acid concentration and purity data for each analyzed sample.
 - **FlowCytometry:** Performs flow cytometry to analyze the characteristics of individual cells or particles in a fluid by passing them through a laser beam, where their size, shape, and fluorescent labels are measured. Expected input materials: One or more cell or particle samples to be analyzed. If action is successful, expected output materials: separate materials for each analyzed sample after the measurement, separate flow cytometry data for each analyzed sample.
 - **Dilute:** Adds a specified amount of solvent to specified samples. Expected input materials: One or more samples to be diluted, solvent. If action is successful, expected output materials: separate materials for each diluted sample, residual solvent (if any).

- 1512 • **Incubate:** Heats and/or mixes the provided samples for a given amount of time at a given
1513 temperature, allowing for a follow up annealing time. Expected input materials: One or
1514 more samples to be incubated. If action is successful, expected output materials: separate
1515 materials for each incubated sample.
- 1516 • **DryingOven:** Dry samples in a drying oven. Expected input materials: One or more
1517 samples to be dried. If action is successful, expected output materials: separate materials
1518 for each dried sample.
- 1519 • **Mix:** Mixes and/or heats the provided samples for a given amount of time at a given rate
1520 and temperature. Expected input materials: One or more samples to be mixed. If action is
1521 successful, expected output materials: separate materials for each mixed sample.
- 1522 • **Combine:** Combines the provided samples into a single sample. Expected input materials:
1523 Two or more samples to be combined. If action is successful, expected output materials:
1524 one combined material containing all input samples.
- 1525 • **Centrifuge:** Spins down the provided samples for a given amount of time at a provided
1526 force or spin rate. Expected input materials: One or more samples to be centrifuged. If
1527 action is successful, expected output materials: separate materials for each centrifuged
1528 sample.
- 1529 • **Degas:** Performs a degassing procedure on the given samples using a specified technique.
1530 Expected input materials: One or more samples to be degassed. If action is successful,
1531 expected output materials: separate materials for each degassed sample.
- 1532 • **Filter:** Passes the provided samples through a given physical filter using a set of optional
1533 different methods. Expected input materials: One or more samples to be filtered. If action is
1534 successful, expected output materials: separate materials for each filtered sample, collected
1535 filtrate materials.
- 1536 • **Autoclave:** Subjects the provided samples or containers to extreme heat and pressure in
1537 order to sterilize. Expected input materials: One or more samples or containers to be
1538 sterilized. If action is successful, expected output materials: separate materials for each
1539 sterilized sample or container.
- 1540 • **Evaporate:** Evaporates solvent from a provided sample under high vacuum at a given
1541 temperature with centrifugation to prevent bumping. Expected input materials: One or
1542 more samples to have solvent evaporated. If action is successful, expected output materials:
1543 separate materials for each sample with reduced solvent volume.
- 1544 • **Lyophilize:** Removes solvents from the provided samples via controlled freezing and
1545 sublimation under high vacuum. Expected input materials: One or more samples to be
1546 lyophilized. If action is successful, expected output materials: separate materials for each
1547 lyophilized sample.
- 1548 • **Aspirate:** Removes the supernatant from a sample by aspiration. Expected input materials:
1549 One or more samples with supernatant to be removed. If action is successful, expected
1550 output materials: separate materials for each sample with supernatant removed, collected
1551 supernatant materials.
- 1552 • **FillToVolume:** Adds sample to the a container until its volume reaches the desired value.
1553 Expected input materials: One or more samples to be filled to volume, solvent for volume
1554 adjustment. If action is successful, expected output materials: separate materials for each
1555 sample adjusted to target volume, residual solvent (if any).
- 1556 • **AcousticLiquidHandling:** Transfers liquid samples with sound waves in nanoliter incre-
1557 ments. Expected input materials: One or more source samples, destination containers. If
1558 action is successful, expected output materials: residual source materials (if any), separate
1559 materials for each destination with transferred amounts.
- 1560 • **AdjustpH:** Adds acid or base titrant to change the pH of the given sample to the desired
1561 value. Expected input materials: One or more samples to have pH adjusted, acid or base
1562 titrant. If action is successful, expected output materials: separate materials for each pH-
1563 adjusted sample, residual titrant (if any).
- 1564
- 1565

- 1566 • **Resuspend:** Dissolve the specified solid samples with some amount of solvent. Can apply
1567 to living and non-living solid samples. Expected input materials: One or more solid sam-
1568 ples to be resuspended, solvent. If action is successful, expected output materials: separate
1569 materials for each resuspended sample, residual solvent (if any).
- 1570 • **ApplyMagnet:** Isolates targets from specified sample via magnetic bead separation, which
1571 uses a magnetic field to separate superparamagnetic particles from suspensions. Expected
1572 input materials: One or more samples containing magnetic beads. If action is successful,
1573 expected output materials: separate materials for each sample with magnetically separated
1574 components.
- 1575 • **Microwave:** Breaks down complex samples via microwave heating and (optional) acid/ox-
1576 idizing agent to fully solubilize sample for subsequent operations, especially ICP-MS. Ex-
1577 pected input materials: One or more samples to be microwaved. If action is successful,
1578 expected output materials: separate materials for each microwaved sample.
- 1579 • **FlashFreeze:** Performs freezing of specified sample objects through immersion of the sam-
1580 ple containers in liquid nitrogen. Expected input materials: One or more samples to be flash
1581 frozen. If action is successful, expected output materials: separate materials for each flash
1582 frozen sample.
- 1583 • **Desiccate:** Dries out solid substances by absorbing water molecules from the samples
1584 through exposing them to a chemical desiccant in a bell jar desiccator under vacuum or
1585 non-vacuum conditions. Expected input materials: One or more samples to be desiccated.
1586 If action is successful, expected output materials: separate materials for each desiccated
1587 sample.
- 1588 • **Grind:** Employs mechanical actions to break particles of solid samples into smaller powder
1589 particles, using a grinding apparatus. Expected input materials: One or more solid samples
1590 to be ground. If action is successful, expected output materials: separate materials for each
1591 ground sample.
- 1592 • **AttemptDNAExtraction:** Series of steps to isolate DNA from a given sample. Also, com-
1593 monly referred to as a mini-prep, midi-prep, etc. Expected input materials: One or more
1594 samples containing DNA to be extracted. If action is successful, expected output materials:
1595 separate materials for each extracted DNA sample.
- 1596 • **CountLiquidParticles:** Measures the number of suspended particles in a liquid colloid
1597 or very fine suspension sample. Expected input materials: One or more liquid samples
1598 with suspended particles to be counted. If action is successful, expected output materials:
1599 separate materials for each analyzed sample after the measurement, separate particle count
1600 data for each analyzed sample.
- 1601 • **CoulterCount:** Measures the number and size distribution of suspended particles (typically
1602 cells) in a liquid colloid or very fine suspension sample. Expected input materials: One or
1603 more samples with suspended particles to be counted. If action is successful, expected out-
1604 put materials: separate materials for each analyzed sample after the measurement, separate
1605 particle count and size distribution data for each analyzed sample.
- 1606 • **MeasureOsmolality:** Measures the concentration of osmotically active species in a solu-
1607 tion. Expected input materials: One or more samples to have osmolality measured. If action
1608 is successful, expected output materials: separate materials for each analyzed sample after
1609 the measurement, separate osmolality data for each analyzed sample.
- 1610 • **MeasureConductivity:** Measures the electrical conductivity of a sample by immersion of
1611 a conductivity probe into the solution. Expected input materials: One or more samples to
1612 have conductivity measured. If action is successful, expected output materials: separate
1613 materials for each analyzed sample after the measurement, separate conductivity data for
1614 each analyzed sample.
- 1615 • **MeasureContactAngle:** Measures the contact angle of a fiber sample with a wetting liquid
1616 using a force tensiometer. Expected input materials: One or more fiber samples, wetting
1617 liquid. If action is successful, expected output materials: separate materials for each ana-
1618 lyzed sample after the measurement, separate contact angle data for each analyzed sample,
1619 residual wetting liquid (if any).

- 1620 • **MeasureDensity**: Measures the density of the given samples using a fixed volume weight
1621 measurement or a density meter. Expected input materials: One or more samples to have
1622 density measured. If action is successful, expected output materials: separate materials
1623 for each analyzed sample after the measurement, separate density data for each analyzed
1624 sample.
- 1625 • **MeasureDissolvedOxygen**: Measures the partial pressure of oxygen in a sample by ap-
1626 plying a constant voltage in a probe confined by an oxygen permeable membrane to detect
1627 oxygen reduction as an electrical signal. Expected input materials: One or more samples to
1628 have dissolved oxygen measured. If action is successful, expected output materials: sepa-
1629 rate materials for each analyzed sample after the measurement, separate dissolved oxygen
1630 data for each analyzed sample.
- 1631 • **MeasurepH**: Measures the pH of the given sample using electrical potential sensors. Ex-
1632 pected input materials: One or more samples to have pH measured. If action is successful,
1633 expected output materials: separate materials for each analyzed sample after the measure-
1634 ment, separate pH data for each analyzed sample.
- 1635 • **MeasureWeight**: Measures the weight of the given samples using an appropriately sized
1636 balance. Expected input materials: One or more samples to have weight measured. If action
1637 is successful, expected output materials: separate materials for each analyzed sample after
1638 the measurement, separate weight data for each analyzed sample.
- 1639 • **MeasureVolume**: Measures the volume of the given samples using ultrasonic measurement
1640 of liquid surface distance and prior parametrization of the surface distance to volume in
1641 the samples container to determine sample volumes. Expected input materials: One or
1642 more samples to have volume measured. If action is successful, expected output materials:
1643 separate materials for each analyzed sample after the measurement, separate volume data
1644 for each analyzed sample.
- 1645 • **MeasureCount**: Measures the number of tablets in a given tablet sample by determining
1646 the average weight of the tablets in the sample and the total mass of the given tablet sample.
1647 Expected input materials: One or more tablet samples to have count measured. If action
1648 is successful, expected output materials: separate materials for each analyzed sample after
1649 the measurement, separate tablet count data for each analyzed sample.
- 1650 • **ImageSample**: Records an image of the given sample either from above or side on for
1651 larger transparent vessels. Expected input materials: One or more samples to be imaged. If
1652 action is successful, expected output materials: separate materials for each imaged sample
1653 after the imaging, separate image data for each imaged sample.
- 1654 • **MeasureSurfaceTension**: Determines the surface tension of a sample by measuring the
1655 forces exerted on a small diameter rod as it is withdrawn from a sample. Expected input
1656 materials: One or more samples to have surface tension measured. If action is successful,
1657 expected output materials: separate materials for each analyzed sample after the measure-
1658 ment, separate surface tension data for each analyzed sample.
- 1659 • **MeasureRefractiveIndex**: Measures the Refractive Index (RI) of the given sample with
1660 refractometer. Expected input materials: One or more samples to have refractive index
1661 measured. If action is successful, expected output materials: separate materials for each
1662 analyzed sample after the measurement, separate refractive index data for each analyzed
1663 sample.
- 1664 • **CyclicVoltammetry**: Characterizes the reduction and oxidation processes of the given
1665 sample using Cyclic Voltammetry (CV). Expected input materials: One or more samples to
1666 be analyzed, reference electrode. If action is successful, expected output materials: sepa-
1667 rate materials for each analyzed sample after the measurement, separate cyclic voltammetry
1668 data for each analyzed sample.
- 1669 • **PrepareReferenceElectrode**: Generates a reference electrode filled with a reference solu-
1670 tion to be used in electrochemical measurements, including Cyclic Voltammetry measure-
1671 ments. Expected input materials: None. If action is successful, expected output materials:
1672 prepared reference electrode.
- 1673 • **VisualInspection**: Monitors the insoluble particles in the given sample while its container
is agitated. Expected input materials: One or more samples to be visually inspected. If ac-

- 1674 tion is successful, expected output materials: separate materials for each inspected sample
1675 after the inspection, separate visual inspection data for each inspected sample.
1676
- 1677 • **MeasureViscosity:** Measures a fluid’s viscosity defined as the resistance to deformation
1678 by assessing the flow rate of the sample when loaded into the viscometer chip. Expected
1679 input materials: One or more fluid samples to have viscosity measured. If action is suc-
1680 cessful, expected output materials: separate materials for each analyzed sample after the
1681 measurement, separate viscosity data for each analyzed sample.
 - 1682 • **DynamicFoamAnalysis:** Characterizes the foamability, stability, drainage process and
1683 structure of liquid-based foams by monitoring foam generation and decay of a sample.
1684 Expected input materials: One or more liquid samples to have foam properties analyzed. If
1685 action is successful, expected output materials: separate materials for each analyzed sample
1686 after the analysis, separate foam analysis data for each analyzed sample.
 - 1687 • **MeasureMeltingPoint:** Measures the melting points of the solid samples using a melting
1688 point apparatus that applies an increasing temperature gradient to melting point capillary
1689 tubes containing a small amount of the input samples. Expected input materials: One or
1690 more solid samples to have melting point measured. If action is successful, expected output
1691 materials: separate materials for each analyzed sample after the measurement, separate
1692 melting point data for each analyzed sample.
 - 1693 • **UseMicroscope:** Performs imaging on provided cellular samples using a bright-field mi-
1694 croscope or a high content imager. Expected input materials: One or more cellular samples
1695 to be imaged. If action is successful, expected output materials: separate materials for each
1696 imaged sample after the imaging, separate microscopy data for each imaged sample.
 - 1697 • **ImageColonies:** Acquires bright-field, absorbance or fluorescence images of the provided
1698 samples containing microbial cells on a solid media plate using a colony handler. Expected
1699 input materials: One or more samples containing microbial colonies to be imaged. If action
1700 is successful, expected output materials: separate materials for each imaged sample after
1701 the imaging, separate colony image data for each imaged sample.
 - 1702 • **FreezeCells:** Lowers the temperature of cell samples under controlled conditions to prepare
1703 cells for long term cryopreservation. Expected input materials: One or more cell samples
1704 to be frozen. If action is successful, expected output materials: separate materials for each
1705 frozen sample.
 - 1706 • **QuantifyColonies:** Measures the microbial cell concentration in the provided samples.
1707 Expected input materials: One or more samples containing microbial colonies to be quanti-
1708 fied. If action is successful, expected output materials: separate materials for each analyzed
1709 sample after the quantification, separate colony count data for each analyzed sample.
 - 1710 • **CellCultureIncubate:** Incubate samples in sterile, controlled conditions. Expected input
1711 materials: One or more cell culture samples to be incubated. If action is successful, ex-
1712 pected output materials: separate materials for each incubated sample.
 - 1713 • **Wait:** Allows the system to idle for a specified amount of time to let materials evolve, react,
1714 or reach equilibrium. Expected input materials: One or more materials to be waited on. If
1715 action is successful, expected output materials: separate materials for each material after
1716 waiting period.
 - 1717 • **Storage:** Store materials under specified environmental conditions for preservation or ap-
1718 propriate handling. Expected input materials: One or more materials to be stored. If action
1719 is successful, expected output materials: separate materials for each stored material.
 - 1720 • **UsePhotospectrometer:** Measure absorbance, transmittance, or fluorescence using photo-
1721 spectrometer. Expected input materials: One or more samples to be analyzed. If action
1722 is successful, expected output materials: separate materials for each analyzed sample after
1723 the measurement, separate spectrophotometric data for each analyzed sample.
 - 1724 • **UsePlateReader:** Measure absorbance, fluorescence, or luminescence using microplate
1725 reader. Expected input materials: One or more samples in microplate format to be analyzed.
1726 If action is successful, expected output materials: separate materials for each analyzed
1727 sample after the measurement, separate plate reader data for each analyzed sample.

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

- **SubmitMaterialsObservationsForReview:** When you believe you have a set of materials and observations that meet the all of the requirements for the goal, submit them for final review.
- **Quit:** If you do not believe you can complete the task with the current materials and observations, quit the process.
- **ManifestMaterials:** Request and add new laboratory materials to the inventory using natural language descriptions. This action has no time cost and does not advance the step count.
- **SeparateWells:** Separates a multi-well plate material into individual well materials or groups based on treatment conditions. Importantly, this is not a laboratory action, but rather a simulation-specific "perception" action to help the simulation agent interact with its available materials by being able to use them as individual samples. The properties of the resulting materials should not change as a result of this action, and no time passes as a result of this action. Rather, they should be properly preserved and simply broken up into individual materials.
- **CellCounter:** Counts the number of cells in a given sample using automated cell counting methods such as hemocytometer, automated cell counter, or flow cytometry. Input material is consumed and is no longer available for use after this method.
- **Discard:** Discards unwanted materials or samples by disposing of them safely according to appropriate waste disposal protocols.