
000 MINIMAL EPISTEMIC CLOSED-LOOP AGENTS FOR SCI-
001 ENTIFIC DISCOVERY: CONSTRAINTS, INFORMATION
002 GAIN, AND DISCRIMINATIVE TESTING IN A TOY GE-
003 NOMICS ENVIRONMENT
004
005
006
007

008 **Anonymous authors**

009 Paper under double-blind review
010
011

012 ABSTRACT
013

014 Closed-loop “AI scientist” systems can generate hypotheses, call tools, and it-
015 erate, but many implicitly optimize for plausibility or predicted success rather
016 than explicitly reducing uncertainty about an underlying mechanism. We study
017 a minimal design for *epistemic* closed-loop agents that (i) enforce feasibility via
018 hard constraints during proposal, (ii) select experiments by Expected Information
019 Gain (EIG), and (iii) accelerate refutation by generating discriminative “Achilles”
020 tests that maximize disagreement among hypotheses. We evaluate only algorithmic
021 behavior in a small, reproducible toy genomics simulator with discrete latent
022 hypotheses and noisy outcomes (not biological validation). Across random seeds,
023 EIG-based selection reduces posterior entropy faster than success-seeking base-
024 lines, and Achilles-style testing reduces experiments-to-refute incorrect leading
025 hypotheses under an operational refutation criterion. Finally, we use a minimal
026 taxonomy for negative outcomes (infeasible, inaccessible, null, execution-failure)
027 to avoid conflating distinct failure semantics during belief updates.
028

029 1 INTRODUCTION
030

031 Closed-loop agents are increasingly capable of drafting hypotheses, invoking tools, and iterating on re-
032 sults (Lu et al., 2024; Bran et al., 2024; Boiko et al., 2023). However, many systems prioritize outputs
033 that appear plausible or likely to succeed under learned heuristics, rather than actions that maximally
034 reduce uncertainty about a hidden mechanism. Bayesian experimental design formalizes epistemic
035 progress via Expected Information Gain (EIG) (Lindley, 1956; Chaloner and Verdinelli, 1995; Rain-
036 forth et al., 2024), and falsification-oriented workflows emphasize actively seeking discriminative
037 tests that separate competing explanations (Liu et al., 2025; Popper, 1959).
038

039 **Scope.** We intentionally restrict claims to a toy genomics simulator to demonstrate algorithmic
040 behavior, not wet-lab capability: (i) hard feasibility constraints, (ii) EIG-based experiment selection,
041 (iii) discriminative “Achilles” tests for rapid refutation, plus a minimal negative-outcome taxonomy.
042

043 **Related work.** Recent “AI scientist” agents demonstrate end-to-end automation across ideation,
044 tool use, and writing (Lu et al., 2024; Bran et al., 2024; Boiko et al., 2023) but do not primarily
045 optimize epistemic uncertainty reduction. Separately, Bayesian experimental design studies EIG
046 objectives and estimators (Lindley, 1956; Chaloner and Verdinelli, 1995; Rainforth et al., 2024),
047 while falsification agents explicitly target refutation behaviors (Liu et al., 2025). Our contribution is a
048 minimal, reproducible closed loop that integrates (constraints + EIG + discriminative testing) and
049 reports simple metrics (entropy reduction and experiments-to-refute) in a controlled simulator.
050

051 2 MINIMAL EPISTEMIC CLOSED-LOOP
052

053 Let $h \in \mathcal{H}$ denote a discrete latent mechanism/hypothesis, $x \in \mathcal{X}$ an experiment, and y an observed
outcome. The agent maintains $p(h \mid \mathcal{D})$ and iteratively selects x to reduce uncertainty.

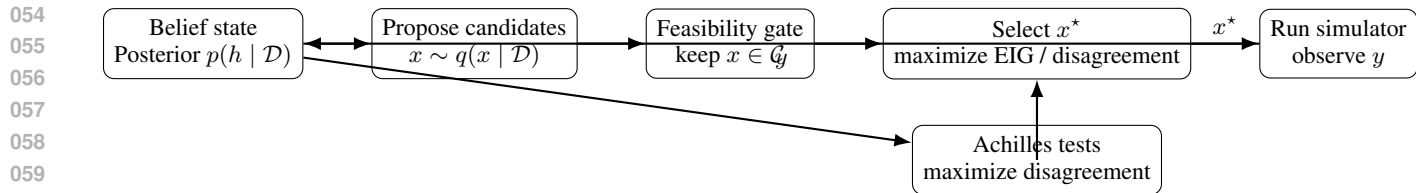


Figure 1: Minimal epistemic closed loop. “Operational refutation” means driving an incorrect leading hypothesis’ posterior below a threshold (Sec. 3).

A: Hard feasibility constraints. Many proposals are invalid due to sequence/protocol/assay constraints. We model feasibility as a hard set \mathcal{C} and do not execute $x \notin \mathcal{C}$. In our toy instantiation, \mathcal{C} includes length bounds, GC bounds, and low-complexity filters. Hard checking can be deterministic or solver-based (De Moura and Bjørner, 2008); here we use deterministic filters.

B: Epistemic planning via EIG. The planner selects experiments maximizing mutual information:

$$x^* = \arg \max_{x \in \mathcal{X}} I(h; y | x, \mathcal{D}). \quad (1)$$

In our discrete toy setting, EIG is computed by enumerating h and marginalizing outcomes.

C: Achilles tests (discriminative testing). Beyond informativeness in expectation, Achilles tests target *disagreement*: propose x that maximally separates predictions under the current leading hypothesis \hat{h} from plausible alternatives, accelerating refutation when \hat{h} is wrong (Liu et al., 2025; Popper, 1959).

3 TOY GENOMICS ENVIRONMENT AND EVALUATION

Environment. We use a synthetic sequence-to-outcome simulator inspired by promoter-like design. Each hypothesis h defines a conditional outcome model $p(y | x, h)$ with noise. The agent begins with a prior over h and updates $p(h | \mathcal{D})$ by Bayes’ rule after observing outcomes.

Baselines. We compare: `random`; `greedy_success` (maximize predicted mean under MAP hypothesis); `uncertainty` (maximize predictive variance); `eig` (Eq. 1); and an `achilles` variant that prioritizes disagreement-based tests.

Metrics. (1) Posterior entropy $H[p(h | \mathcal{D})]$ vs. step. (2) Experiments-to-refute: steps until an incorrect leading hypothesis satisfies $p(\hat{h} | \mathcal{D}) < \epsilon$.

3.1 RESULTS

EIG-based selection reduces posterior entropy faster than `random` and `greedy_success`, consistent with EIG optimizing expected posterior change (Rainforth et al., 2024). Achilles-style selection reduces experiments-to-refute by prioritizing queries that maximize predictive separation across hypotheses.

4 NEGATIVE OUTCOMES AND FAILURE SEMANTICS

Negative outcomes are not monolithic: conflating them can induce incorrect belief updates. We use a minimal four-way taxonomy: **(i) infeasible** (violates \mathcal{C}), **(ii) inaccessible** (assay limitation / unobservable), **(iii) null** (no signal under noise), **(iv) execution failure** (the experiment did not execute correctly. Execution failures are treated as epistemically null events and do not contribute evidence for or against any hypothesis.) This supports principled updates and motivates structured logging of negative outcomes (Fanelli, 2012; Ioannidis, 2005).

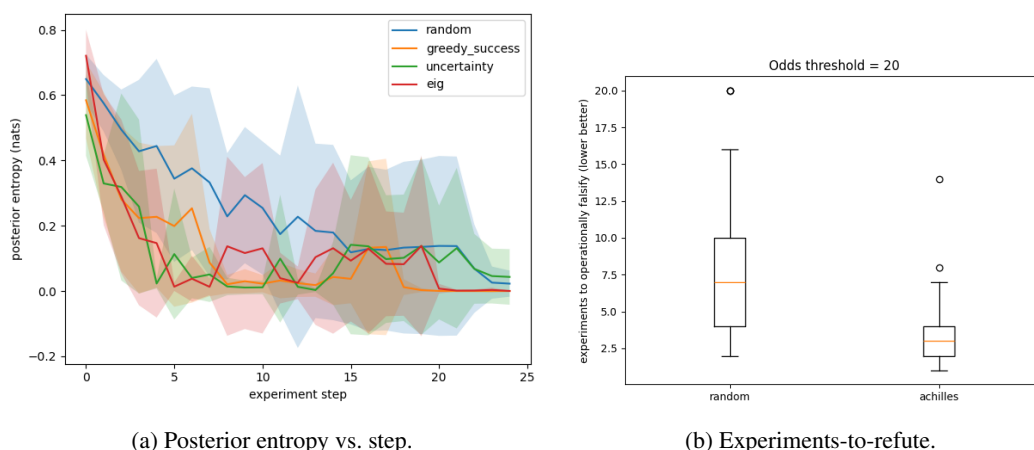


Figure 2: Toy environment results.

5 LIMITATIONS AND CONCLUSION

This paper demonstrates algorithmic behavior in a small simulator and does not claim biological validation or autonomous wet-lab capability. Within this scope, combining feasibility constraints, EIG-based selection, and discriminative testing yields more epistemic progress per experiment. Next steps include scaling to established closed-loop benchmarks (Gandhi et al., 2025; Jansen et al., 2024) and testing whether failure semantics improve calibration and decision quality.

REFERENCES

- Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., and Ha, D. The AI Scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- Bran, A.-M., Cox, S., Schilter, O., Baldassari, C., White, A. D., and Schwaller, P. ChemCrow: Augmenting large-language models with chemistry tools. *Nature Machine Intelligence*, 6:525–535, 2024.
- Boiko, D. A., MacKnight, R., Kline, B., and Gomes, G. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, 2023.
- Lindley, D. V. On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, 27(4):986–1005, 1956.
- Chaloner, K. and Verdinelli, I. Bayesian experimental design: A review. *Statistical Science*, 10(3):273–304, 1995.
- Rainforth, T., Foster, A., Ivanova, D. R., and Bickford Smith, F. Modern Bayesian experimental design. *Statistical Science*, 39(1):100–127, 2024.
- Liu, Z., Liu, K., Zhu, Y., Lei, X., Yang, Z., Zhang, Z., Li, P., and Liu, Y. BABY-AIGS: An LLM-agent framework with explicit falsification for AI-generated science. *arXiv preprint*, 2025.
- Popper, K. *The Logic of Scientific Discovery*. Hutchinson, 1959.
- De Moura, L. and Bjørner, N. Z3: An efficient SMT solver. In *TACAS*, pp. 337–340, 2008.
- Fanelli, D. Negative results are disappearing from most disciplines and countries. *Scientometrics*, 90(3):891–904, 2012.
- Ioannidis, J. P. Why most published research findings are false. *PLoS Medicine*, 2(8):e124, 2005.

162 Gandhi, K., Li, M. Y., Goodyear, L., Li, L., Bhaskar, A., Zaman, M., and Goodman, N. D. Box-
163 ingGym: Benchmarking progress in automated experimental design and model discovery. *arXiv*
164 *preprint*, 2025.

165
166 Jansen, P. A., Côté, M.-A., Khot, T., Bransom, E., Dalvi, B., Majumder, B. P., Tafjord, O., and Clark,
167 P. DiscoveryWorld: A virtual environment for developing and evaluating automated scientific
168 discovery agents. *arXiv preprint*, 2024.

169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215