

# Unlearning as Ablation: Toward a Falsifiable Benchmark for Generative Scientific Discovery

Robert Yang

Independent Researcher

San Jose, CA 95129

bobyang9@alumni.stanford.edu

## Abstract

Bold claims about AI’s role in science—from “AGI will cure all diseases” to promises of radically accelerated discovery—raise a central epistemic question: do large language models (LLMs) truly *generate* new knowledge, or do they merely remix memorized fragments? We propose **unlearning-as-ablation** as a falsifiable probe of constructive scientific discovery. The idea is to systematically remove a target result together with its *forget-closure* (supporting lemmas, paraphrases, and multi-hop entailments) and then evaluate whether the model can re-derive the result from only permitted axioms and tools. Success would indicate generative capability beyond recall; failure would expose current limits. Unlike prevailing motivations for unlearning—privacy, copyright, or safety—our framing repositions it as an *epistemic probe* for AI-for-Science. We outline a minimal pilot in mathematics and algorithms to illustrate feasibility, and sketch how the same approach could later be extended to domains such as physics or chemistry. This is a position paper: our contribution is conceptual and methodological, not empirical. We aim to stimulate discussion on how principled ablation tests could help distinguish models that reconstruct knowledge from those that merely retrieve it, and how such probes might guide the next generation of AI-for-Science benchmarks.

## 1 Introduction

Recent breakthroughs in foundation models have fueled bold claims—from predictions that “AGI will cure all diseases” to assertions that scientific progress will soon accelerate far beyond historical rates. These visions reflect real excitement, but they obscure a fundamental epistemic question: **do large language models (LLMs) genuinely generate new knowledge, or do they merely remix what was already present in their training data?**

This distinction matters deeply for AI-for-Science. Without a falsifiable test of constructive knowledge generation, claims of “discovery” remain philosophically ambiguous and scientifically ungrounded. If AI systems are to be trusted as collaborators in science, we must know whether they can *derive* new results from principles, rather than retrieve or interpolate memorized fragments.

We propose a new perspective: **unlearning-as-ablation**. The idea is straightforward. Select a scientific result  $T$  (e.g., a theorem or algorithm), identify its entire *forget-closure*  $\mathcal{F}(T)$ —all lemmas, paraphrases, aliases, and multi-hop entailments that lead to  $T$ —and perform strong unlearning over  $\mathcal{F}(T)$ . Afterward, provide the model only with permitted axioms and tools, and test whether it can re-derive  $T$  in a verifiable form. Success constitutes positive evidence of constructive generation, whereas failure or leakage exposes the boundaries of current capabilities.

This framing departs from prevailing motivations for unlearning. Surveys emphasize privacy, copyright, and safety as primary rationales [Xu et al., 2023, 2024], with evaluation focused on removal

fidelity rather than generative ability. Recent work highlights the difficulty of faithfully removing multi-hop or entangled knowledge [Choi et al., 2024, Wang et al., 2025, Shah et al., 2025], while other studies show that forgotten content can often be “relearned” through small finetunes or prompting [Shah et al., 2025]. In the safety and compliance setting, these phenomena are treated as risks. In our setting, they define the frontier: as unlearning methods improve in addressing leakage and robustness, the resulting ablations become more faithful, and the corresponding rediscovery benchmarks more stringent. In this way, progress in unlearning directly strengthens our ability to test whether models are capable of constructive scientific generation.

By reframing unlearning as an *experimental probe*, we aim to bridge AI-for-Science and safety communities. The result is a concrete, falsifiable methodology for testing the limits of LLMs: whether they are capable of genuine discovery, or whether their advances remain bounded by retrieval and interpolation. As a position paper, our contribution is primarily conceptual: we propose a methodological framework and outline pilot domains, leaving systematic empirical validation to future work.

## 2 Background: Unlearning Today

The study of unlearning in machine learning and large language models (LLMs) has grown rapidly in recent years, motivated largely by *external constraints* such as law, safety, or ethics rather than by epistemic goals. We briefly review the dominant rationales, common methodologies, and key evaluation challenges.

### 2.1 Motivations for Unlearning

Three primary motivations recur across surveys and frameworks:

**(1) Privacy and compliance.** Regulations such as the General Data Protection Regulation (GDPR) enshrine a “right to be forgotten,” requiring that models support the removal of sensitive or personally identifiable data. Surveys on digital forgetting in LLMs emphasize compliance with privacy law as a central driver of research in this area [Xu et al., 2024].

**(2) Copyright and intellectual property.** LLMs trained on large web scrapes may inadvertently memorize copyrighted text, code, or images. Several works argue that machine unlearning is necessary to respect intellectual property claims and to support takedown requests from rights-holders [Karamolegkou et al., 2023, Dou et al., 2025, Xu et al., 2024, Ren et al., 2025, Yao et al., 2024].

**(3) Safety and dual-use knowledge.** A third line of work focuses on removing *hazardous* content: for example, step-by-step instructions for synthesizing explosives or pathogens. Recent benchmarks such as WMDP [Li et al., 2024] evaluate whether unlearning can reduce dual-use risks while maintaining general utility.

### 2.2 Methodological Approaches

Most unlearning methods adapt techniques from model editing or fine-tuning. Examples include:

- **Gradient-ascent or anti-training:** adjusting model parameters to maximize loss on target examples, thereby forgetting them.
- **Representation-level interventions:** e.g., Amnesic Probing [Elazar et al., 2021] removes specific linguistic features from hidden states.
- **Retrieval suppression:** steering methods that block particular outputs without removing underlying representations.

While diverse, these approaches generally aim at *removal fidelity*: ensuring that specific facts or behaviors no longer appear in model outputs.

### 2.3 Evaluation Challenges

Evaluation is a persistent bottleneck. Several recent studies emphasize that:

- **Entangled knowledge is difficult to erase.** Multi-hop unlearning benchmarks show that even if intermediate nodes are removed, models can often reconstruct targets via alternative reasoning chains [Choi et al., 2024, Wang et al., 2025, Shah et al., 2025].
- **Suppression vs. removal.** SoK papers stress the importance of distinguishing true parameter-level removal from surface-level suppression, where models appear to forget but can be prompted to recall [Ren et al., 2025].
- **Relearning and robustness.** Empirical work demonstrates that forgotten content can often be “jogged” back into use with minimal finetuning or prompting [Lee et al., 2025].

## 2.4 Gap for AI-for-Science

Notably, none of the above rationales frame unlearning as a tool for *scientific epistemology*. Unlearning has been motivated by compliance and safety, not by the question of whether a model can *reconstruct* forgotten knowledge from first principles. This gap opens an opportunity: by treating unlearning as *ablation*, we can design falsifiable experiments to probe whether LLMs possess constructive generative capabilities, a perspective particularly urgent for AI-for-Science. Moreover, the progress of unlearning research directly determines the strength of such benchmarks: the more thorough and faithful the unlearning, the harder the rediscovery task becomes, and the more reliable the test of whether models can generate knowledge rather than recall it.

## 3 Proposal: Unlearning-as-Ablation

We propose to repurpose unlearning from its conventional role in privacy or safety into an *experimental ablation method* for probing constructive knowledge generation. The central idea is to remove not only a target result  $T$ , but also all of the *supporting knowledge that directly enables it*, and then ask the model to re-derive  $T$  from only axioms and tools that remain accessible. If the model succeeds under these conditions, we gain falsifiable evidence that it is not merely retrieving memorized fragments but genuinely generating knowledge.

### 3.1 Defining the Forget-Closure

The first step is to formally define the **forget-closure**  $\mathcal{F}(T)$  of a target  $T$ . This closure includes:

- All direct statements of  $T$  (canonical forms, proofs, code).
- Paraphrases and rephrasings that preserve semantic equivalence.
- Intermediate lemmas or building blocks that entail  $T$ .
- Multi-hop reasoning chains where  $T$  can be reconstructed indirectly [Choi et al., 2024, Wang et al., 2025, Shah et al., 2025].
- Same-answer sets where multiple formulations yield equivalent outputs.

By removing the entire  $\mathcal{F}(T)$ , we close off not only surface forms but also indirect reasoning paths that would otherwise allow reconstruction through entanglement.

### 3.2 Performing Strong Unlearning

The second step is to apply **removal-oriented unlearning** across  $\mathcal{F}(T)$ . Unlike suppression methods that steer generation away from target outputs, removal aims to eliminate relevant information from the parameterization itself. Candidate techniques include gradient-ascent unlearning, targeted finetuning, or optimization-based methods evaluated in recent surveys [Ren et al., 2025]. To confirm removal, we propose adopting multi-faceted audits:

- Leakage checks on paraphrase, multi-hop, and same-answer sets.
- Counterfactual activation probes (inspired by Amnesic Probing) to test whether  $T$ -related features still reside in hidden states [Elazar et al., 2021].
- Robustness tests against “jogging” attacks, where small finetunes or prompting can restore forgotten knowledge [Lee et al., 2025].

These checks ensure that the unlearning process produces a genuine epistemic blank slate with respect to  $\mathcal{F}(T)$ .

### 3.3 Re-Derivation as a Falsifiable Test

Finally, we design a **re-derivation trial**. After unlearning, the model is provided with:

1. A set of axioms, primitives, or base tools that are *not* part of  $\mathcal{F}(T)$ .
2. A prompt or environment that permits constructive reasoning (e.g., a proof assistant or a test-driven code synthesis framework).

The task is to derive  $T$  in a form that can be verified by an external oracle: for example, a formal proof accepted by Lean or Isabelle, or a program passing a hidden test suite. Importantly, success is only counted if  $T$  is re-derived *without leakage from  $\mathcal{F}(T)$* .

This yields a falsifiable criterion: if the model can re-derive  $T$  despite rigorous unlearning of all prerequisite paths, we have positive evidence for constructive generation. If it cannot, or if leakage audits reveal dependence on residual memory, then the claim of “scientific discovery” remains unsubstantiated.

### 3.4 Why This Matters

This approach connects progress in unlearning directly to progress in measuring scientific discovery. In the safety and compliance literature, challenges such as entanglement, multi-hop reasoning, and relearning are treated as failure modes because they undermine removal fidelity [Dou et al., 2025, Choi et al., 2024, Wang et al., 2025, Shah et al., 2025]. In our framing, they set the difficulty of the benchmark: the more effectively unlearning methods address these challenges, the more thoroughly the target knowledge is ablated, and the more demanding the rediscovery task becomes. Thus, advances in unlearning translate into sharper tests of whether LLMs truly possess constructive generative capability. Rather than turning flaws into benefits, we highlight that solving these long-standing problems in unlearning is what enables rigorous epistemic evaluation in AI-for-Science.

## 4 Minimal Pilot Study

While the long-term vision is to apply unlearning-as-ablation to scientific hypotheses in physics, chemistry, or biology, we propose beginning with domains where **verification is automatic and unambiguous**. This allows us to isolate the epistemic question—can a model *re-derive* knowledge once its closure has been forgotten?—without relying on subjective human judgment.

### 4.1 Mathematics: Formal Proofs

Mathematics provides an ideal testbed because results can be verified by proof assistants such as *Lean* or *Isabelle*. A minimal pilot could proceed as follows:

1. Select a mid-tier theorem (e.g., in number theory or combinatorics) that has a clear dependency structure.
2. Construct its forget-closure  $\mathcal{F}(T)$ , including canonical statements, paraphrased variants, and prerequisite lemmas.
3. Apply strong unlearning over  $\mathcal{F}(T)$ .
4. Task the model with re-proving  $T$  using only base axioms and allowed rules of inference.

Success is defined as producing a proof accepted by the proof assistant. Failure or leakage (e.g., shortcut recall of a forgotten lemma) falsifies the claim of rediscovery.

### 4.2 Algorithms: Verified Implementations

Algorithms provide another tractable domain, where correctness can be checked against hidden test suites. For example:

1. Forget the Knuth–Morris–Pratt (KMP) string matching algorithm, along with all prerequisite explanations, code templates, and paraphrases.
2. After unlearning, ask the model to derive an efficient string-matching procedure from first principles (e.g., reasoning about prefix functions).
3. Validate correctness using adversarial test cases and runtime complexity checks.

As in mathematics, the evaluation is binary: either the model reconstructs a working implementation, or it does not.

### 4.3 Evaluation Metrics

To assess the outcome of such pilots, we propose three classes of metrics:

- **Success rate.** Fraction of trials where the model re-derives  $T$  in a verifiable form (proof acceptance, program passes test suite).
- **Leakage audits.** Performance on paraphrase, multi-hop, and same-answer sets drawn from  $\mathcal{F}(T)$ , ensuring the model is not recalling forgotten material [Choi et al., 2024, Wang et al., 2025, Shah et al., 2025].
- **Utility retention.** Accuracy on unrelated benchmarks (e.g., a subset of MMLU) to confirm that unlearning did not degrade general capability [Ren et al., 2025, Yao et al., 2024].

### 4.4 Why a Minimal Pilot is Valuable

Even small-scale pilots can decisively answer whether LLMs exhibit generative capability under ablation. If a model successfully re-derives a theorem or algorithm after strong unlearning of its closure, we obtain falsifiable evidence that it constructs knowledge rather than merely retrieving it. Conversely, if models fail under such controlled conditions, this highlights a concrete epistemic limit of current systems. Either outcome offers high-value insight for AI-for-Science, where claims of accelerated discovery remain both enticing and contested.

## 5 Implications for AI-for-Science

The proposed unlearning-as-ablation framework has direct consequences for how we understand the promise and limits of AI-for-Science.

### 5.1 Epistemic Clarity in Scientific Discovery

The central value of this approach is that it provides a *falsifiable test* of discovery. Today, when an LLM proposes a hypothesis, proves a theorem, or writes an algorithm, it remains unclear whether this is a product of genuine reasoning or of subtle retrieval from training data. By first *removing* all accessible pathways to a result and then testing for *re-derivation*, we create a clean epistemic separation: success implies constructive generation, while failure implies dependence on stored fragments. This reframing allows the AI-for-Science community to move beyond speculation about “discovery” and instead ground claims in falsifiable evidence.

### 5.2 Turning Failure Modes into Probes

Unlearning research has traditionally cast entanglement, multi-hop reasoning, and relearning as obstacles [Choi et al., 2024, Wang et al., 2025, Shah et al., 2025]. In our setting, these challenges become useful stress tests. If a model cannot succeed once closure paths are blocked, it indicates that the relevant knowledge was never truly generative. If it can succeed, it demonstrates robustness and constructive capacity. Either way, phenomena previously treated as evaluation headaches become diagnostic instruments for probing the depth of model reasoning.

### 5.3 Broader AI-for-Science Roadmap

Although we highlight mathematics and algorithms as tractable pilot domains, the methodology generalizes. In physics, one could remove an established equation and test whether the model can

re-derive it from fundamental laws. In chemistry, one could unlearn a well-known synthesis route and test whether the model can rediscover it from reaction rules. In biology, one could unlearn a canonical protein interaction and test for re-derivation from structural principles. These extensions would demand careful closure construction and domain-specific verification, but they illustrate how the same ablation logic scales to real scientific practice.

#### 5.4 Redefining the Boundary of AI Progress

Finally, this framework speaks directly to the theme of this workshop: the reach and limits of AI in scientific discovery. If unlearning-as-ablation pilots reveal that models can re-derive knowledge under strong ablation, this strengthens the case that AI can generate truly novel insights. If they reveal consistent failures, it delineates a boundary condition: LLMs may accelerate retrieval, interpolation, and synthesis, but fall short of independent knowledge generation. In both outcomes, the methodology provides a principled way to map the contours of what AI can and cannot do for science.

#### 5.5 Toward the Next Major Benchmark

A final implication is that unlearning-as-ablation offers a clear path toward the next generation of benchmarks for AI progress. Just as ImageNet catalyzed advances in computer vision by providing a well-defined task on which algorithms could be compared [Deng et al., 2009], a benchmark grounded in constructive re-derivation after unlearning could serve as a lodestar for AI-for-Science. Existing evaluations of knowledge regurgitation and short-form reasoning are increasingly saturated—as highlighted by works such as Humanity’s Last Exam [Phan et al., 2025]—suggesting that the next frontier must measure whether models can move beyond retrieval and interpolation to genuine discovery. We believe that such an “unlearning-as-ablation” benchmark could become a distinguishing test of model strength, separating systems that can merely recall from those that can constructively generate new scientific knowledge.

Importantly, the strength of such a benchmark is coupled to the progress of unlearning research itself. As unlearning methods become more faithful and thorough, the corresponding benchmarks become more stringent: rediscovery requires deeper reasoning, and successful re-derivation provides stronger evidence of constructive capability. In this way, advances in unlearning directly drive advances in our ability to measure—and eventually to achieve—genuine AI scientific discovery.

### 6 Conclusion

We have proposed *unlearning-as-ablation* as a new lens on large language models, reframing unlearning from a tool of compliance and safety into a falsifiable probe of scientific discovery. By systematically removing a target result and its forget-closure, and then testing whether the model can re-derive the result from permitted axioms and tools, we obtain an experimental method to separate retrieval from constructive generation. This approach directly addresses one of the most pressing open questions in AI-for-Science: can AI systems truly generate new knowledge? Even minimal pilots in mathematics or algorithms provide decisive evidence either way, while extensions to physics, chemistry, and biology can delineate the boundaries of future AI scientific progress. Whether the outcome is success or failure, unlearning-as-ablation offers the community a principled framework to move beyond speculation and anchor claims of discovery in falsifiable tests.

### References

- Minseok Choi, ChaeHun Park, Dohyun Lee, and Jaegul Choo. Breaking chains: Unraveling the links in multi-hop knowledge unlearning. 10 2024. doi: 10.48550/arXiv.2410.13274.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, and Eric Wong. Avoiding copyright infringement via large language model unlearning. pages 5176–5200, 01 2025. doi: 10.18653/v1/2025.findings-naacl.288.

- Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic probing: Behavioral explanation with amnesic counterfactuals. *Transactions of the Association for Computational Linguistics*, 9:160–175, 03 2021. ISSN 2307-387X. doi: 10.1162/tac1\_a\_00359. URL [https://doi.org/10.1162/tac1\\_a\\_00359](https://doi.org/10.1162/tac1_a_00359).
- Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. Copyright violations and large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7403–7412, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.458. URL <https://aclanthology.org/2023.emnlp-main.458/>.
- Bruce W. Lee, Addie Foote, Alex Infanger, Leni Shor, Harish Kamath, Jacob Goldman-Wetzler, Bryce Woodworth, Alex Cloud, and Alexander Matt Turner. Distillation robustifies unlearning, 2025. URL <https://arxiv.org/abs/2506.06278>.
- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Gabriel Mukobi, Nathan Helm-Burger, Rassin Lababidi, Lennart Justen, Andrew Bo Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xiaoyuan Zhu, Rishub Tamirisa, Bhruhu Bharathi, Ariel Herbert-Voss, Cort B Breuer, Andy Zou, Mantas Mazeika, Zifan Wang, Palash Oswal, Weiran Lin, Adam Alfred Hunt, Justin Tienken-Harder, Kevin Y. Shih, Kemper Talley, John Guan, Ian Steneker, David Campbell, Brad Jokubaitis, Steven Basart, Stephen Fitz, Ponnurangam Kumaraguru, Kallol Krishna Karmakar, Uday Tupakula, Vijay Varadharajan, Yan Shoshitaishvili, Jimmy Ba, Kevin M. Esvelt, Alexandr Wang, and Dan Hendrycks. The WMDP benchmark: Measuring and reducing malicious use with unlearning. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 28525–28550. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/li24bc.html>.
- Long Phan, Andrew Gatti, Zihao Han, et al. Humanity’s last exam. 2025. URL <https://arxiv.org/abs/2501.14249>.
- Jie Ren, Yue Xing, Yingqian Cui, Charu C Aggarwal, and Hui Lui. Sok: Machine unlearning for large language models. 06 2025. doi: arXiv.2506.09227.
- Raj Sanjay Shah, Jing Huang, Keerthiram Murugesan, Nathalie Baracaldo, and Diyi Yang. The unlearning mirage: A dynamic framework for evaluating LLM unlearning. In *Second Conference on Language Modeling*, 2025. URL <https://openreview.net/forum?id=exW2SFJK4H>.
- Changsheng Wang, Chongyu Fan, Yihua Zhang, Jinghan Jia, Dennis Wei, Parikshit Ram, Nathalie Baracaldo, and Sijia Liu. Reasoning model unlearning: Forgetting traces, not just answers, while preserving reasoning skills. 06 2025. doi: arXiv.2506.12963.
- Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S. Yu. Machine unlearning: A survey. *ACM Comput. Surv.*, 56(1), August 2023. ISSN 0360-0300. doi: 10.1145/3603620. URL <https://doi.org/10.1145/3603620>.
- Jie Xu, Zihan Wu, Cong Wang, and Xiaohua Jia. Machine unlearning: Solutions and challenges. *IEEE Transactions on Emerging Topics in Computational Intelligence*, PP:1–19, 06 2024. doi: 10.1109/TETCI.2024.3379240.
- Yuanshun Yao, Xiaojun Xu, and YangLiu. Large language model unlearning. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 105425–105475. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/be52acf6bccf4a8c0a90fe2f5cfcead3-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/be52acf6bccf4a8c0a90fe2f5cfcead3-Paper-Conference.pdf).