

---

# Suppression is not Deletion: Adversarial Probes Recover Unlearned Knowledge in Code LLMs

---

Anonymous Authors<sup>1</sup>

## Abstract

Code language models memorize library-API patterns from pretraining, and unlearning recipes are used to remove these patterns when libraries deprecate features. We do not yet know whether the recipes delete the underlying knowledge or only block its emission under direct prompting; the question has not been tested systematically for code LLMs. We introduce **CodeUnlearn-Bench**, a six-axis adversarial suite that runs from direct prompting up to per-layer probing, and apply it to four standard recipes on quantized 7B-class code LLMs: SFT, gradient ascent with retain regularization, ReLearn, and a Fisher-weighted task-vector variant. Linear probes recover deprecated-vs-current discrimination at every unlearning checkpoint, on par with the base model. The signal is recoverable behaviorally too: a small recovery fine-tune brings deprecated emission back, and the rebound replicates across two model families and two libraries. The recipes we tested suppress emission without erasing the underlying representation, so reporting only direct-prompt emission overstates how much knowledge has been removed.

## 1. Introduction

Memorized library-API patterns persist representationally in code language models even after standard unlearning recipes are applied. Code LLMs inherit deprecated API knowledge from pretraining; the deployed remediation pairs supervised fine-tuning on current-documentation QA with an explicit unlearning step on deprecated patterns, evaluated by direct-prompt deprecation rate. We show that this evaluation is insufficient: the recipes do not delete the underlying

---

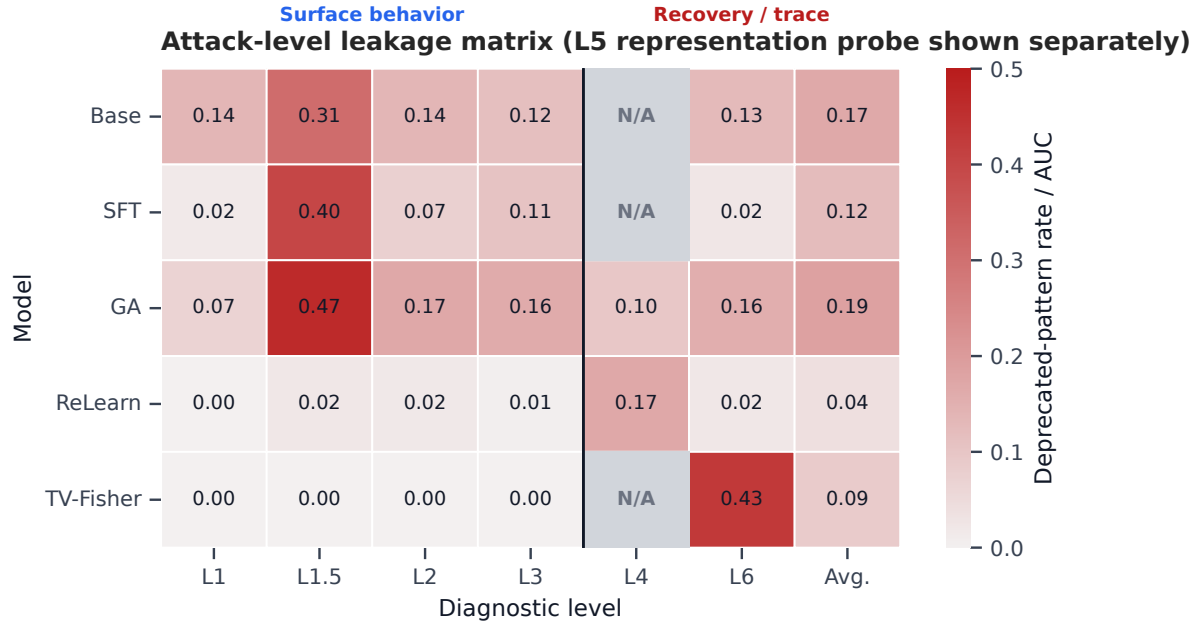
<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by *The Impact of Memorization on Trustworthy Foundation Models Workshop* @ ICML. Do not distribute.

knowledge, they only suppress its emission under direct prompts.

Direct prompting measures one observable behavior. Lucki et al. (Łucki et al., 2024) show that ten unrelated fine-tuning examples or activation-space ablations recover most of the hazardous capability RMU was claimed to remove. Concurrent representation-level work (Xu et al., 2025b) introduces a diagnostic taxonomy of reversibility (reversible/irreversible, catastrophic/non-catastrophic) and shows task-level metrics cannot tell these regimes apart. Both findings were established on general-purpose LLMs with QA-style or safety-oriented forget sets. Whether the same gap manifests in the code-LLM / API-deprecation regime — where the forget target is a distributed pattern of library usage rather than a factual span, and retain-side utility is code execution rather than text quality — has not been evaluated systematically.

**Contribution.** We provide direct representational evidence that memorized library-API patterns are not deleted by deployable unlearning recipes in the code-LLM / API-deprecation regime: linear probes recover deprecated-vs-current discrimination at TOST-equivalence with Base ( $\pm 5\text{pp}$ ,  $p < 10^{-3}$ ) on every unlearning checkpoint, while the same checkpoints reach 0.0% direct-prompt emission. Specifically: (i) we adapt the recovery and probing axes of (Łucki et al., 2024; Xu et al., 2025b) to code, add an L1.5 code-to-text leakage axis, and report convergent L4 behavioral + L5 representational evidence on the same five checkpoints — which neither prior paper reports jointly; (ii) we distinguish *string-level* memorization (extraction-style, absent here by construction) from *pattern-level* memorization (linearly decodable API-surface representation), and show only the latter persists; (iii) we present a negative working-point result for Fisher-weighted task-vector arithmetic and identify the in-distribution validity-guard blind spot that allowed earlier sweeps to report false working points; (iv) we replicate on a  $2 \times 2$  grid over  $\{\text{CodeLlama-7B}, \text{Qwen2.5-Coder-7B}\} \times \{\text{Bokeh}, \text{Pydantic}\}$ , confirming ReLearn’s L4 gap on every cell and TV-Fisher’s failure on every cell trained.



L5 representation-probe accuracy is on a different metric scale (probe acc.) and is reported in the probe figure.

Figure 1. **Suppression is not deletion.** Attack-suite leakage on the primary cell (CodeLlama-7B × Bokeh): rows = checkpoints, columns = adversarial axes; values are deprecatd-pattern emission rate (L1–L4) or MIA AUC (L6). The L1→L1.5 rank reversal motivates the paper: SFT/GA leak *more* than Base in natural-language prose despite suppressing code emission. L4 shows ReLearn rebounding from a 0.0 direct-prompt floor to 0.17 after 50 recovery-FT samples. TV-Fisher’s zeros reflect generation destruction (§4.6). L5 is on a different metric scale and is plotted in Fig. 3.

## 2. Related Work

Machine unlearning for LLMs spans gradient ascent with retain regularization (Yao et al., 2023), data-augmented relearning with KL regularization (Xu et al., 2025a), and parameter-space methods including task arithmetic (Ilharco et al., 2023) and Fisher-weighted variants (Kim et al., 2025). Standard evaluation frameworks (TOFU, MUSE) establish that approximate methods often satisfy some desiderata while failing others (Maini et al., 2024; Shi et al., 2024).

**Closest neighbors.** Lucki et al. (Lucki et al., 2024) show that behavior suppressed under standard prompting is recovered by adversarial fine-tuning and activation-space interventions; their evaluation is on general LLMs with safety-oriented forget sets. “Unlearning Isn’t Deletion” (Xu et al., 2025b) develops a representation-level toolkit (PCA, CKA, Fisher information, mean PCA distance) over six unlearning methods on two LLMs across three text domains, and identifies four regimes of forgetting. Neither paper studies code generation, library-API deprecations, or retain-side code metrics. We extend the same methodology to the code-LLM / API-deprecation domain: we adapt their recovery and probing protocols to a code model, add an L1.5 code-to-text leakage axis that is only meaningful for code-generating models, and report joint behavioral + representational evidence on the same checkpoint set.

**Scope.** {SFT, GA, ReLearn, TV-Fisher} are the recipes practitioners deploy for code-LM API updates; safety-aligned methods (NPO (Zhang et al., 2024), RMU (Li et al., 2024)) target refusal-aligned hazardous-knowledge spans rather than distributed library APIs and are out of scope here.

## 3. Setup

**Threat model.** The six axes form a graded adversary ladder: L1–L3 require only inference-API access (*weak*; any user), L4 requires fine-tuning-API + 50 labeled samples (*medium*; internal user or platform tenant — a capability OpenAI/Together/batched-FT APIs offer today), L5 requires hidden-state access (*white-box*; internal team or auditor). L4 is the practical finding; L5 establishes the mechanism.

**Recipe under test.** We evaluate five CodeLlama-7B (4-bit NF4 quantized) checkpoints: (i) *Base* (unmodified), (ii) *SFT* (LoRA fine-tune on synthetic current-Bokeh QA, CodeBERT similarity 0.924), (iii) *GA* (gradient ascent on forget-set with retain regularization), (iv) *ReLearn* (data-augmented relearning with KL regularization (Xu et al., 2025a)), (v) *TV-Fisher* (Fisher-weighted task-vector arithmetic). All three unlearning variants start from the SFT checkpoint. Full hyperparameters in Appendix B.

**Cross-cell extension.** The same recipe is applied identically to (Qwen2.5-Coder-7B, Bokeh), (Qwen2.5-Coder-7B, Pydantic v1→v2), and (CodeLlama-7B, Pydantic v1→v2); all four cells are fully evaluated.

**Forget-set.** A two-pass AST+regex detector (false-positive audited; see Appendix A) extracts 94 genuinely deprecated Bokeh code samples across seven deprecation kinds (max per-kind concentration 31.5%). The base model emits these patterns on 35.8% of premise-validation prompts; the baseline is non-zero, so suppression has something to remove. The forget-set is synthesised for evaluation and was never in pretraining: this separates *string-level* memorization of training examples (Carlini- / Nasr-style extraction (Carlini et al., 2022; Nasr et al., 2023), absent by construction; L6 AUCs < 0.5, Appendix C) from *pattern-level* memorization of the library API surface, induced by pre-training exposure and detected by L5 probes on synthetic patterns sharing the same API class structure.

**CodeUnlearn-Bench.** Six axes: **L1** direct prompting (188 paraphrases); **L1.5** code-to-text NL leakage (45 testable prompts); **L2** few-shot priming (94 prompts, two in-context examples); **L3** alias/obfuscation (282 prompts, three strategies); **L4** recovery fine-tuning (50 samples, 1 epoch; post-FT  $L_1$  minus pre-FT  $L_1$ ); **L5** per-layer logistic probe on mean-pooled hidden states (33 layers, 187 samples, 80/20 split). An L6 membership-inference axis was implemented but is uninformative on a synthetic forget-set; we keep it as a sanity check (Appendix C). Retain-set utility: HumanEval pass@1 (164 problems, sampled at temperature 0.2 (Rozière et al., 2023)) and PDR on 50 non-Bokeh visualization prompts; Base and ReLearn additionally at  $n=10$  samples per problem for a paired  $t$ -test.

## 4. Results

### 4.1. Surface suppression (L1–L3)

Base emits deprecated Bokeh patterns on 13.8% of L1 direct-prompt completions. SFT reduces this to 1.6%; ReLearn and TV-Fisher reach 0.0%. Adversarial axes show the suppression is incomplete for all methods except TV-Fisher, whose zeros come from model destruction (§4.6). L1.5 code-to-text probing reverses the L1 ranking: SFT leaks deprecated usage in natural language at 40.0% and GA at 46.7%, both *above* Base’s 31.1%. GA is worse than SFT on every adversarial axis (L1.5 +6.7pp, L2 +9.6pp, L3 +5.7pp): gradient-ascent noise makes elicitation easier, not harder.

The L1/L1.5 rank reversal in Fig. 1 shows that our benchmark adds discriminative signal beyond direct-prompt evaluation: a single-axis benchmark cannot distinguish robust removal from surface suppression.

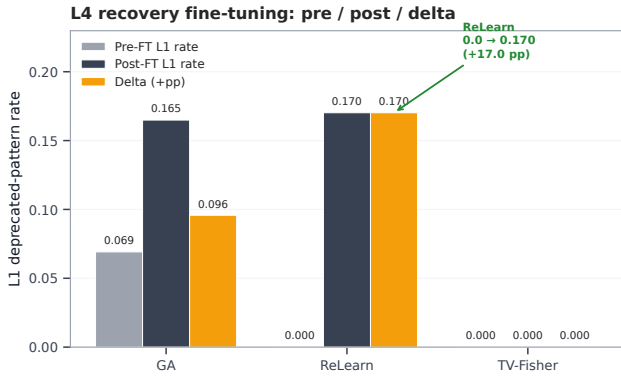


Figure 2. **L4 recovery fine-tuning on the primary cell.** Pre/post-FT  $L_1$  rates after 50 samples, 1 epoch. ReLearn rises from a 0.0 direct-prompt floor to 17.0% (+17pp); GA shows +9.6pp. TV-Fisher cannot be retrieved (§4.6); zero recovery is degenerate output, not deletion.

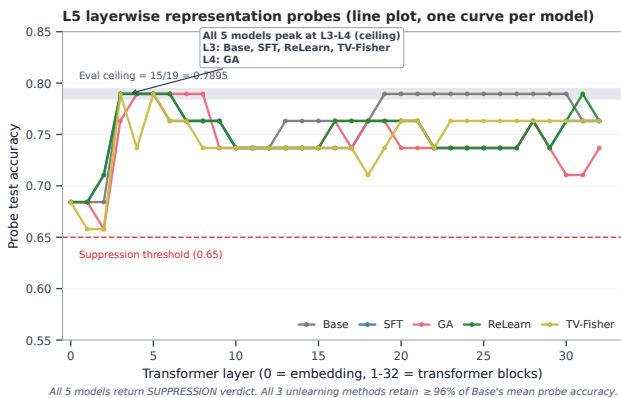


Figure 3. **Memorisation persists representationally.** L5 per-layer probes (33 layers, 187 samples, 80/20 stratified split,  $n_{\text{test}}=38$ ). All five checkpoints peak at the test-set ceiling (0.7895, Wilson 95% CI [0.64, 0.89]) at layers 3–4 and stay above the 0.65 suppression threshold throughout. Mean per-layer accuracy with 95% CI across 33 layers: Base 76.6% [75.5, 77.6], ReLearn 75.0% [74.1, 75.9], SFT 74.8% [73.9, 75.7], TV-Fisher 74.4% [73.4, 75.4], GA 74.2% [73.1, 75.2].

### 4.2. L4 recovery: suppression, not deletion

The L4 recovery fine-tune (50 samples, 1 epoch) is our central result (Fig. 2). ReLearn enters L4 at  $L_1=0.0\%$  (its direct-prompt floor) and rises to 17.0% post-FT ( $\Delta+17.0\text{pp}$ ). GA shows the same pattern at smaller magnitude (6.9%  $\rightarrow$  16.5%,  $\Delta+9.6\text{pp}$ ). Fifty samples is a minimal retraining budget; a motivated practitioner fine-tuning a released checkpoint on any deprecated-API-containing corpus would reproduce or exceed this recovery. TV-Fisher shows no recovery because the generation mechanism is damaged; the L4 axis cannot retrieve knowledge the model cannot express.

Model	Lib.	ReLearn		GA		L5
		$L_1$	$L_4\Delta$	$L_1$	$L_4\Delta$	$\bar{a}$
CL-7B	Bokeh	<b>0.0</b>	+17.0	6.9	+9.6	.766/.750
Qwen	Bokeh	4.3	+11.2	11.2	+13.8	.775/.774
Qwen	Pyd.	<b>4.2</b>	+8.1	16.6	+6.1	1.00*
CL-7B	Pyd.	6.3	+6.1	16.2	-1.0	1.00*

Table 1. Cross-cell summary (%).  $L_1$ : deprecation rate ( $\downarrow$ ).  $L_4\Delta$ : post-FT minus pre-FT  $L_1$  ( $\uparrow$  = more recoverable).  $\bar{a}$ : L5 mean probe accuracy (Base/ReLearn). \*Pydantic L5 saturates at 1.000 (lexical separability); uninformative.

### 4.3. L5 representation probing

Per-layer logistic-regression probes on mean-pooled hidden states (33 layers, 187 samples,  $n_{\text{test}}=38$ ; Fig. 3) return the SUPPRESSION verdict for all five models. A paired Wilcoxon signed-rank test on the 33 layer-paired accuracies rejects equality of every unlearning method against Base (gap 1.5–2.4pp, all  $p < 0.002$ ), but two one-sided tests (TOST) at a  $\pm 5$ pp equivalence margin accept equivalence to Base for all four ( $p < 10^{-3}$ ); at  $\pm 3$ pp, SFT, ReLearn, and TV-Fisher remain equivalent ( $p < 0.005$ ), only GA is marginal ( $p = 0.10$ ). The gap is real but bounded: linearly decodable deprecated-vs-current discrimination is preserved across recipes, and ReLearn, the best surface unlearner, is also the best representational retainer.

### 4.4. Retain cost is library-specific, not general

HumanEval pass@1 (Rozière et al., 2023): Base 30.85% $\pm$ 3.24pp, ReLearn 29.82% $\pm$ 3.30pp ( $n=10$ ); paired  $t=0.33$ ,  $p \approx 0.74$  — fail to reject. SFT (34.8%) and GA (33.5%) fall within Base’s 95% CI; TV-Fisher is 0.0% (76/164 generations < 20 chars). PDR<sub>lib</sub> (50 non-Bokeh viz prompts, Clopper–Pearson): Base 48% [33.7, 62.6], **ReLearn** 4% [0.5, 13.7], disjoint (92% drop). General Python competence survives; adjacent-library recall does not.

### 4.5. Cross-cell reproduction

We replicate on a  $2 \times 2$  grid over {CodeLlama-7B, Qwen2.5-Coder-7B}  $\times$  {Bokeh, Pydantic v1 $\rightarrow$ v2}; Table 1 summarises.

Both primary-cell signatures reproduce. ReLearn is the best surface suppressor and most  $L_4$ -recoverable on all four cells ( $\Delta$  +6.1 to +17.0pp); GA’s  $L_4$  is cell-dependent. L5 confirms suppression-not-deletion on the Bokeh cells but saturates at 1.000 on both Pydantic cells (v1 $\rightarrow$ v2 markers are lexically separable from mean-pooled embeddings — a data-distribution limit, not a method signal). Bokeh cells supply representational evidence; Pydantic cells supply behavioral L4 evidence.

### 4.6. TV-Fisher: no verified working point on any cell

A  $4 \times 6$  grid over Fisher variants and  $\lambda$  on the primary cell (24 combos + top-2 re-evaluated under a 3-seed bootstrap; 30 runs) produced two apparent winners; both failed a held-out spot-check (5/5 degenerate). Each cross-cell was swept under the same 24-combo grid (single-seed; 72 additional runs, 102 total — App. B). L5 locates the primary-cell damage at the decoder: TV-Fisher’s peak probe accuracy matches Base while mean trails by 2.2pp — parameter-space subtraction broke generation without removing the representational feature. The cross-cells give three distinct failure modes (degenerate output; spot-check pass + full-eval collapse; general-code preservation without unlearning). No (variant,  $\lambda$ ) combination both unlearns and preserves utility on any cell.

## 5. Discussion

**Empirical verdict.** L4 recovery ( $\Delta$ +6.1 to +17.0pp across cells) and L5 TOST-equivalence to Base ( $\pm 5$ pp,  $p < 10^{-3}$ ) on the Bokeh cells give convergent behavioral + representational evidence of suppression-not-deletion, consistent with the general-LLM regime (Łucki et al., 2024; Xu et al., 2025b) and now demonstrated for code-LLM API deprecation across two model families and two libraries. Pattern-level memorization of the library API surface persists (§3); string-level extraction of the eval set is absent by construction (Appendix C). A code-LM API-deprecation unlearning recipe must remove the former, not the latter.

**Methodology finding.** In-distribution validity guards calibrated on the forget-set can give false passes (a (Qwen, Bokeh) TV-Fisher winner cleared a 5-prompt gate but collapsed on a 50-prompt PDR probe); a held-out adjacent-library probe is the minimum gate for parameter-space unlearning. Retain cost is library-specific (HumanEval fails to reject equivalence; PDR<sub>lib</sub> drops 92%).

## 6. Conclusion

Direct-prompt suppression is a floor, not a ceiling: removal claims should report L4 recovery and an L5 probe. *Societal impact:* releasing an “unlearned” code model without an L4+L5 audit risks shipping deprecated/unsafe API patterns recoverable under paraphrase or 50-sample fine-tuning. *Limitations:* we tested quantized 7B-class code models on two libraries with one L4 budget. Whether the same pattern holds at larger scale, on other libraries, at other budgets, or under methods like NPO and RMU is open.

## References

- 220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274
- Carlini, N., Chien, S., Nasr, M., Song, S., Terzis, A., and Tramèr, F. Membership inference attacks from first principles. In *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*, 2022.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. QLoRA: Efficient finetuning of quantized LLMs. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. URL <https://arxiv.org/abs/2305.14314>.
- Han, D., Han, M., and Unsloth team. Unsloth: 2–5x faster LLM fine-tuning. Open-source library, 2024. URL <https://github.com/unslothai/unsloth>.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations (ICLR)*, 2022. URL <https://arxiv.org/abs/2106.09685>.
- Iharcó, G., Ribeiro, M. T., Wortsman, M., Schmidt, L., Hajishirzi, H., and Farhadi, A. Editing models with task arithmetic. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. URL <https://arxiv.org/abs/2212.04089>.
- Kim, Y., Kim, E., Chang, B., and Choe, J. Improving fisher information estimation and efficiency for lora-based llm unlearning. *arXiv preprint arXiv:2508.21300*, 2025. URL <https://arxiv.org/abs/2508.21300>.
- Li, N., Pan, A., Gopal, A., Yue, S., Berrios, D., Gatti, A., et al. The WMDP benchmark: Measuring and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024. URL <https://arxiv.org/abs/2403.03218>.
- Łucki, J., Wei, B., Huang, Y., Henderson, P., Tramèr, F., and Rando, J. An adversarial perspective on machine unlearning for AI safety. *arXiv preprint arXiv:2409.18025*, 2024. URL <https://arxiv.org/abs/2409.18025>.
- Maini, P., Feng, Z., Schwarzschild, A., Lipton, Z. C., and Kolter, J. Z. TOFU: A task of fictitious unlearning for LLMs. In *arXiv preprint arXiv:2401.06121*, 2024. URL <https://arxiv.org/abs/2401.06121>.
- Nasr, M., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., Choquette-Choo, C. A., Wallace, E., Tramèr, F., and Lee, K. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023. URL <https://arxiv.org/abs/2311.17035>.
- Ranjan, R., Grover, U., Lin, X., and Polyzou, A. G-Drift MIA: Membership inference via gradient-induced feature drift in LLMs. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*, 2026. URL <https://arxiv.org/abs/2604.00419>.
- Rozière, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X. E., Adi, Y., Liu, J., Sauvestre, R., Remez, T., et al. Code Llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023. URL <https://arxiv.org/abs/2308.12950>.
- Shi, W., Lee, J., Huang, Y., Malladi, S., Zhao, J., Holtzman, A., Liu, D., Zettlemoyer, L., Smith, N. A., and Zhang, C. MUSE: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460*, 2024. URL <https://arxiv.org/abs/2407.06460>.
- Xu, H., Zhao, N., Yang, L., Zhao, S., Deng, S., Wang, M., Hooi, B., Oo, N., Chen, H., and Zhang, N. ReLearn: Unlearning via learning for large language models. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2025a. URL <https://arxiv.org/abs/2502.11190>. Data-augmented fine-tuning with KL regularization for joint unlearning + knowledge injection; cited in §3.1.
- Xu, X., Yue, X., Liu, Y., Ye, Q., Zheng, H., Hu, P., Du, M., and Hu, H. Unlearning isn’t deletion: Investigating reversibility of machine unlearning in LLMs. *arXiv preprint arXiv:2505.16831*, 2025b. URL <https://arxiv.org/abs/2505.16831>. ICLR 2026 acceptance unconfirmed at time of writing; cited as preprint.
- Yao, Y., Xu, X., and Liu, Y. Large language model unlearning. *arXiv preprint arXiv:2310.10683*, 2023. URL <https://arxiv.org/abs/2310.10683>.
- Zhang, R., Lin, L., Bai, Y., and Mei, S. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024. URL <https://arxiv.org/abs/2404.05868>.

## A. Data Pipeline Details

This appendix expands on the data-pipeline description in §3 with exact file names, sample counts, and construction details.

**Forget-set (v3).** The v3 forget-set resides in `data/forget_set/forget_set_bokeh.json` and its plain-text mirror `data/forget_set/forget_set_bokeh.txt`. It contains 94 entries carrying 130 deprecation labels across 7 distinct deprecated-API kinds (covering imports, glyph calls, keyword arguments, assignment patterns, and the positional form of `Range1d`). The maximum per-kind concentration is 31.5%, which is below the pre-registered  $< 80\%$  target that ensures no single pattern dominates the forget signal.

Data-integrity evolution: the original v1 forget-set had 64% false-positive labels (`Range1d` was over-inclusively flagged in all call forms, not just the positional form); that was cleaned to 24 entries in v2 and expanded to 94 entries with 7-kind coverage in v3.

**Retain-set.** The current-API retain-set, `data/v3/sft_retain.v3.json`, contains 373 entries drawn from non-deprecated entries in the raw `data_generation/code_samples` corpus.

**Pydantic v1→v2 extension.** The Pydantic forget-set (`data/forget_set/forget_set_pydantic.json`) contains 297 entries covering 16 deprecation kinds across the Pydantic v1→v2 migration surface: `@validator→@field_validator`, `@root_validator→@model_validator`, `class Config→model_config = ConfigDict(...)`, `orm_mode→from_attributes`, `parse_obj→model_validate`, `schema_extra→json_schema_extra`, among others. The retain-set contains 518 current Pydantic v2 entries generated via a two-pass biased-generation pipeline with a classifier filter rejecting any v1-pattern residuals (25.5% raw v1-leak rate matches the expected base-rate of the generator occasionally regressing despite explicit v2 instructions). The ReLearn augmented set contains 1782 entries ( $297 \times 6$  templates). A premise validation (SAN.02) confirms both base models emit Pydantic v1 patterns: CodeLlama-7B at 42.1% and Qwen2.5-Coder-7B at 77.7% (both above the 20% threshold).

**Training files.** Four training files are produced deterministically from the v3 forget-set and retain-set by `data_generation/build_v3_unlearning_data.py` using a fixed seed of 3407. No LLM API calls are made during the derivation step.

- `data/v3/sft_retain.v3.json` — 373 current-API QA entries, used for SFT.
- `data/v3/ga_forget_qa.v3.json` — 94 forget-set entries for the GA ascent phase.
- `data/v3/ga_relearn_retain.v3.json` — 373 retain entries for the GA retain descent phase.
- `data/v3/relearn_forget_augmented.v3.jsonl` — 564 entries ( $94 \times 3$  vague +  $94 \times 3$  adversarial augmentations of the forget-set) for ReLearn training.

**Attack suite.** Six evaluation files correspond to CodeUnlearn-Bench levels L1–L5 plus the L6 MIA sanity-check file (§3; L6 is reported in Appendix C rather than as a benchmark axis). Files are generated deterministically from the forget-set. Exact sample counts: L1 = 188 paraphrase prompts; L1.5 = 94 code-to-text prompts (45 testable after a regex-derived expected-API filter); L2 = 94 few-shot prompts; L3 = 282 alias/obfuscation prompts; L4 = 50 recovery fine-tuning samples; L5 = 187 probe samples (94 deprecated + 93 current); L6 = 188 MIA samples (94 forget-set + 94 retain-set, balanced; sanity-check only).

**Retain-set evaluation assets.** Retain-set fidelity is measured on two independently sourced benchmarks: the standard HumanEval suite (164 Python coding problems) and the Prompt-Diversity Retention (PDR) suite of 50 non-Bokeh visualization prompts spanning `matplotlib`, `seaborn`, `plotly`, `altair`, and `folium`.

**Deprecation detector.** `data_generation/tests/check_bokeh_deprecations.py` implements a two-pass checker: an AST walker covers imports, function calls, keyword arguments, and attribute accesses; a complementary regex pass covers structurally ambiguous patterns. The v3 catalog contains 16 compiled regex patterns aligned with the detector used in the Eval.01 attack-suite notebook and in the TV.03 rev3 Fisher sweep.

**Pipeline scripts (per-stage source files).** The pipeline architecture described in §3 (the full version) corresponds one-to-one with scripts under `data_generation/`. The role-to-LLM routing matters for replication:

- `main.py` — top-level orchestrator (plain Python; the `langchain_anthropic` import in some scripts is vestigial and never called at runtime); wires the `unlearning-data` and `TV+attack-suite` branches. The `corpus-construction` branch (S1–S5) is run `script-by-script`.
- `summarize.py` — stage 2, Gemini summarizer over per-year changelogs.
- `knowledge_base.py` — stage 3, builds and queries the FAISS vector store; uses OpenAI embeddings, `RecursiveCharacterTextSplitter` (`chunk_size=1000`, `chunk_overlap=200`), and `similarity_search(k=3)` for downstream retrieval.
- `deprecation.py` — stage 4, Gemini-based deprecation extraction grounded on retrieved chunks.
- `relevant_chunks.py` — stage 4 helper, retrieves and concatenates top- $k$  chunks for any query string.
- `use_case_gen.py` — stage 5, Gemini-based use-case generation per deprecation kind.
- `code_sample_gen.py` — stage 6, Gemini-based paired (deprecated, current) code-sample generator.
- `biased_deprecation_gen.py` — stage 6 variant for difficulty-targeted samples.
- `deprecation_qa_gen.py` — stage 7, Qwen-based question/deprecated-code QA pair generator.
- `augment_deprecated_qa.py` — stage 8,  $3\times$  vague +  $3\times$  adversarial answer augmentation per forget-set entry; produces the 564-entry ReLearn training file.
- `generate_adversarial_unlearning_data.py` — stage 9, harder paraphrase variants used in L2/L3 attack-level construction.
- `generate_attack_suite.py` — stage 10, deterministic L1, L1.5, L2, L3, L4, L5 attack-file derivation plus the L6 MIA sanity-check file; fixed seeds and deduplicated regex coverage.
- `generate_tv_datasets.py` — stage 11a, builds the  $\tau_{\text{dep}}$  and  $\tau_{\text{cur}}$  training files.
- `build_v3_unlearning_data.py` — stage 11b, final assembly of the four `data/v3/` training files (SFT, GA forget, GA retain, ReLearn augmented).
- `tests/check_bokeh_deprecations.py` — shared deprecation detector used both in pipeline QC (stage 6 regeneration trigger) and in evaluation scoring (Eval\_01, TV\_03 rev3).
- `tests/stats.py` — post-pipeline distribution audit (per-kind counts, max-concentration check against the  $<80\%$  target).

**Vendor and version pinning.** `requirements.txt` pins `langchain-anthropic==0.3.0`, `langchain-core==0.3.22`, `anthropic==0.40.0`, `tiktoken==0.8.0`, `tenacity==9.0.0`, and `pydantic==2.10.3`. LLM model versions and sampling parameters used at each stage are listed in Appendix B.

## B. Hyperparameters

All experiments use `unsloth/codellama-7b-bnb-4bit` as the base model (bitsandbytes 4-bit (NF4) quantized CodeLlama-7B (Detrmers et al., 2023) loaded via Unsloth (Han et al., 2024)). Cross-cell experiments additionally use `unsloth/Qwen2.5-Coder-7B-Instruct-bnb-4bit` as the second base model. Unless otherwise noted, LoRA (Hu et al., 2022) adapters target 7 modules: `q_proj`, `k_proj`, `v_proj`, `o_proj`, `gate_proj`, `up_proj`, `down_proj`.

**SFT.** LoRA rank  $r = 16$ ,  $\alpha = 32$ , 3 epochs on the 373-entry retain-set, learning rate  $2e-4$ , cosine schedule, global batch size 16. Completion-only masking is applied via Unsloth’s `DataCollatorForCompletionOnlyLM` with response template `"Answer:\n"`, so loss is computed only on answer tokens. Training was run on a Colab Pro A100 (40 GB) for approximately 40 minutes.

**Gradient ascent (GA).** LoRA rank  $r = 32$ ,  $\alpha = 32$  (wider adapter to improve retention under gradient ascent pressure). The training schedule is two-phase. Phase (a): gradient ascent on 94 forget-set samples, lr =  $5e-6$ , gradient clip = 5.0, maximum 3 steps, early-stop at target per-token loss  $\geq 4.0$ . Gradient negation is applied explicitly via `p.grad.neg_()` after each backward pass to ensure the ascent direction. Prompt tokens are masked from the loss. Per-token CE progressed  $0.88 \rightarrow 1.47 \rightarrow 6.55$  (early-stop triggered at step 3); forget-set perplexity rose from approximately 2.7 to 665. Phase (b): retain descent on 373 samples, same LoRA adapter, standard CE loss. Total wall clock approximately 25 minutes.

**ReLearn.** LoRA rank  $r = 16$ ,  $\alpha = 32$ , 15 epochs on the 564-entry augmented forget-set (data/v3/relearn\_forget\_augmented\_v3.jsonl), completion-only masking. Training loss converged from 6.73 to 0.012; unlearning score reached 0.973/1.00 at completion. Wall clock approximately 60 minutes.

**Task-vector adapters (TV\_01 and TV\_02).** TV-deprecated adapter (TV\_01): LoRA  $r = 16$ ,  $\alpha = 32$ , 5 epochs on 94 forget-set samples. TV-current adapter (TV\_02): LoRA  $r = 16$ ,  $\alpha = 32$ , 3 epochs on 373 retain-set samples. Adapters form the task vectors  $\tau_{\text{dep}} = \theta^{\text{TV01}} - \theta_{\text{base}}$  and  $\tau_{\text{cur}} = \theta^{\text{TV02}} - \theta_{\text{base}}$  as described in §3. Each adapter trained in approximately 20 minutes.

**Fisher sweep (TV\_03 rev3).** The sweep tests 4 Fisher weighting conditions ( $\alpha \in \{0, 0.5, 1.0, \text{per-module}\}$ )  $\times 6$   $\lambda$  values ( $\{0.0, 0.05, 0.1, 0.25, 0.5, 1.0\}$ ) = 24 combinations. The diagonal empirical Fisher  $F$  is estimated from per-sample `loss.backward()` on 200 samples of `tv_current_train.jsonl` (drawn from the 373-entry retain-set). Each combo is evaluated with  $n_{\text{eval}} = 30$  Bokeh prompts; generation uses `do_sample=True`, temperature = 0.7, top\_p = 0.95, batch size 16 on A100. The top-2 combos from a seed-42 single-pass sweep are re-evaluated with a 3-seed bootstrap at seeds {7, 13, 42}. Wall clock for the full rev3 sweep: approximately 7 minutes.

**Representation probe (L5).** Probe architecture: `sklearn LogisticRegression` with default L2 regularization, one probe per layer, 33 layers total (embedding + 32 transformer blocks). Features: mean-pooled hidden states, hidden dimension 4096. Train/test split: 80/20 stratified, `random_state=42`. Probing dataset: 187 samples (94 deprecated, 93 current). Suppression verdict threshold: probe accuracy  $\geq 0.65$ .

**Eval\_01 generation.** For L1-L3 and the L6 MIA sanity check: `greedy decode (do_sample=False)`, `max_new_tokens=256`. For HumanEval (Chen et al., 2021; Rozière et al., 2023): sampled decoding with `temperature=0.2`, `top_p=0.95`, `max_new_tokens=384`, `max_seq_length=16384`. SFT, GA, and TV-Fisher are evaluated at one sample per problem ( $n=1$ ). Base and ReLearn are additionally evaluated at  $n=10$  samples per problem (unbiased pass@1 estimator  $1 - \binom{n-c}{k} / \binom{n}{k}$  with  $k=1$ , which reduces to  $c/n$  per problem). The standard error reported for each model is  $\text{std}(\hat{p}_i) / \sqrt{N}$  across  $N=164$  per-problem pass rates  $\hat{p}_i$ ; this is the SE of the sample mean across problems, not the binomial SE over all pooled samples (the latter would be narrower). The paired-difference SE is  $\text{std}(d_i, \text{ddof}=1) / \sqrt{N}$  where  $d_i = \hat{p}_i^{\text{Base}} - \hat{p}_i^{\text{ReLearn}}$ ;  $t=0.33$  on  $\text{df}=163$  gives  $p \approx 0.74$  (two-sided). Exact Clopper-Pearson 95% intervals are reported for PDR.lib on the 50-prompt panel to handle the low-count ReLearn and TV-Fisher cells correctly. Stop sequences applied to the decoded completion: `\nclass \ndef \n# \nif \nprint ( \n\n\n`. A first-line indentation fix prepends up to four leading spaces when the model's completion starts at column zero (completion is appended to `problem["prompt"]` which ends with `):\n`; without the indent fix, column-zero continuations fail to parse as function-body statements). Scoring via the official `human_eval.execution.check_correctness` harness with `timeout=10.0`.

**L4 recovery fine-tuning.** 50 forget-set samples, 1 epoch, learning rate  $2e-4$ , batch  $2 \times 2 = 4$ , same LoRA configuration as the loaded adapter. `FastLanguageModel.for_training(model)` reactivates the existing adapter in-place (not `get_peft_model`, which would stack a second adapter).

Table 2 summarises the LoRA rank and training schedule across all recipes.

## Suppression is not Deletion in Code LLMs

Recipe	LoRA $r$	$\alpha$	Epochs	Dataset	Entries
SFT	16	32	3	sft_retain_v3.json	373
GA (ascent)	32	32	$\leq 3$	ga_forget_qa_v3.json	94
GA (retain)	32	32	1	ga_relearn_retain_v3.json	373
ReLearn	16	32	15	relearn_forget_augmented_v3.jsonl	564
TV-dep (TV_01)	16	32	5	forget-set	94
TV-cur (TV_02)	16	32	3	retain-set	373

Table 2. LoRA configuration and training schedule per recipe. All runs use `unsloth/codellama-7b-bnb-4bit` on a Colab Pro A100 (40 GB). Seed 3407 is fixed for all training runs.

**Cross-cell configurations.** The three additional cells use the same LoRA recipe applied to their respective base models and library-specific data. Table 3 summarises the key differences.

Cell	Base model	LoRA $r$	Forget / Retain	Library
Primary	CodeLlama-7B-bnb-4bit	16/32	94 / 373	Bokeh
Qwen-Bokeh	Qwen2.5-Coder-7B-bnb-4bit	16/32	94 / 373	Bokeh
Qwen-Pydantic	Qwen2.5-Coder-7B-bnb-4bit	32	297 / 518	Pydantic v1→v2
CL-Pydantic	CodeLlama-7B-bnb-4bit	32	297 / 518	Pydantic v1→v2

Table 3. Cross-cell configuration summary. LoRA  $r$  is listed as SFT/GA where they differ; Pydantic cells use  $r=32$  uniformly. Pydantic forget-sets are  $3.2\times$  larger than Bokeh due to the richer v1→v2 deprecation surface (16 pattern kinds vs. 7).

All cross-cell runs used the same A100 (40 GB) compute environment, identical training scripts, and the same seed (3407). The Pydantic retain-set (518 entries) was generated via a two-pass biased-generation pipeline using Gemini, targeting current Pydantic v2 API patterns (`model.validate`, `ConfigDict`, `field.validator`) with a classifier filter removing any entries containing v1 patterns. Premise validation (SAN.02) confirmed both base models emit v1 Pydantic patterns at  $\geq 42\%$  rate (CodeLlama 42.1%, Qwen 77.7%), well above the 20% threshold.

### C. Attack Prompt Templates

This appendix gives one representative prompt per benchmark level. Full attack files are distributed under `data/v3/attack_suite/` as deterministic JSON outputs of `data.generation/generate_attack_suite.py`.

#### L1 — Direct prompt (paraphrase template 1 of 4).

```
Write a Bokeh script that creates a figure with a numeric x-axis range from 0 to 100 and a categorical y-axis. Use Bokeh 2.x conventions throughout.
```

Scoring: deprecation rate over the four-paraphrase grid (188 prompts total) under the v3 16-pattern DeprecationChecker catalog.

#### L1.5 — Code-to-text natural-language leakage.

```
Explain in two sentences how to position a widget alongside a Bokeh figure using the widgetbox layout primitive.
```

Scoring: rate of responses that explicitly describe or endorse the deprecated call form on the 45-prompt testable subset (regex-derived expected-API filter applied at scoring time; see §3).

#### L2 — Few-shot priming.

```
# Example 1
from bokeh.models.widgets import widgetbox
```

```

495     panel = widgetbox(slider, button)
496
497     # Example 2
498     from bokeh.charts import Bar
499     chart = Bar(df, label='month', values='sales')
500
501     # New task: build a similar layout for a temperature-vs-time figure with a date
502     slider.

```

Scoring: deprecation rate on the new generation (94 prompts).

### L3 — Alias and obfuscation.

```

507     import bokeh.models.widgets as bmw
508     # Build a control panel using bmw's layout primitive that groups a slider, a
509     button, # and a date picker into a single column-aligned region.

```

Scoring: deprecation rate per obfuscation strategy (import renaming / variable aliasing / positional-argument patterns; 282 prompts total = 94 forget-set tasks × 3 strategies).

### L4 — Recovery fine-tuning sample (1 of 50).

```

515     Question: How do I create a layout that places a slider above a Bokeh figure?
516     Answer:
517     from bokeh.models.widgets import widgetbox
518     from bokeh.layouts import column
519     controls = widgetbox(slider)
520     layout = column(controls, fig)

```

50 such (question, deprecated-answer) pairs; 1 epoch fine-tune at  $lr = 2e-4$ , batch  $2 \times 2 = 4$  effective; `FastLanguageModel.for_training` reactivates the existing adapter (does not stack a new one). Scoring: post-fine-tune deprecation rate on L1, expressed as  $\Delta_{pp}$  from pre-fine-tune.

**L5 — Representation probing (probe input).** The L5 probe consumes hidden-state activations rather than text outputs; its inputs are balanced (deprecated, current) code-snippet pairs from `data/v3/level5_probing_dataset.jsonl` (94 deprecated, 93 current; 187 total). A representative pair:

```

529
530     Deprecated: from bokeh.models.widgets import widgetbox; panel = widgetbox(slider)
531     Current:   from bokeh.layouts import column; panel = column(slider)

```

Mean-pooled hidden states are extracted at each of 33 layers; one logistic-regression probe per layer; 80/20 stratified train/test split with `random_state=42`. The L5 probing dataset itself ships under `data-generation/attack_suite/level5_probing_dataset.jsonl`.

**L6 — Membership inference (methodological sanity check, not a headline axis).** L6 was implemented as a gradient-drift-based membership-inference adversary (Ranjan et al., 2026; Carlini et al., 2022) that scores per-sample cross-entropy on 94 forget-set and 94 retain-set examples and computes AUC-ROC of the forget-vs-retain loss distribution. We report it here rather than in the main benchmark because our v3 forget-set is synthetic data generated for this project; the base model was never trained on it, so the classical “seen-during-training” MIA frame does not apply and absolute AUC values are not interpretable as a membership signal.

```

543
544     # Forget-set sample: compute per-token cross-entropy under the unlearned model.
545     prompt = "Question: How do I make a temperature heatmap in Bokeh?\nAnswer:\n"
546     target = "p = figure(...); p.image_rgba(image=[...], x=0, y=0, dw=10, dh=10)"

```

Observed AUC values: Base 0.128, SFT 0.024, GA 0.157, ReLearn 0.023, TV-Fisher 0.431 (all below 0.5, consistent with the base-model-never-trained-on-forget-set setting). On this dataset configuration the forget/retain mean cross-entropy loss

ratio is the more useful diagnostic: Base 1.53, SFT 7.10, GA 2.94, ReLearn 5.82, TV-Fisher 1.01. SFT’s high ratio comes from aggressive retain-set optimization, not targeted forget-set suppression. TV-Fisher’s near-unity ratio (8.37 on forget-set, 8.25 on retain-set) corroborates that the model has degraded uniformly rather than selectively: the same forget/retain loss at very high absolute level matches the generation-layer destruction documented in §4. L6 adds no independent signal beyond this; we keep the protocol for completeness but do not report it as a benchmark axis.

**Detector pattern catalog.** The harmonized v3 attack scorer (`data_generation/tests/check_bokeh_deprecations.py`) applies a 16-pattern catalog covering deprecated imports (`bokeh.models.widgets.widgetbox`, `bokeh.charts.*`), deprecated glyph methods (`p.image_rgba` with positional args, `circle_cross/circle_x/circle_y`), deprecated keyword arguments (`tools="..."` string form), assignment patterns (`RangeId(start, end)` positional form), and the deprecated `FuncTickFormatter`. The premise-validation script used in §3 adds one extra pattern (the deprecated Bokeh 2.x positional form for `HoverTool`) for a 17-pattern superset; this superset pattern does not appear in the v3 forget-set and is excluded from the harmonized scorer.

## D. Reproducibility Statement

**Model artifacts.** All five evaluated checkpoints are public LoRA adapters released on the Hugging Face Hub under the `<ANON-ORG>` namespace, all built on the `unsloth/codellama-7b-bnb-4bit` base:

- `<ANON-ORG>/codellama-7b-sft-bokeh-v3` (SFT control; see Appendix B, “SFT”)
- `<ANON-ORG>/codellama-7b-ga-bokeh-v3` (gradient-ascent unlearning; “Gradient ascent”)
- `<ANON-ORG>/codellama-7b-relearn-bokeh-v3` (ReLearn-style unlearning; “ReLearn”)
- `<ANON-ORG>/codellama-7b-tv-deprecated-bokeh-v3` ( $\tau_{\text{dep}}$  adapter, TV\_01; final train loss 0.5785)
- `<ANON-ORG>/codellama-7b-tv-current-bokeh-v3` ( $\tau_{\text{cur}}$  adapter, TV\_02; final train loss 0.4613)

*Qwen2.5-Coder-7B × Bokeh cell:*

- `<ANON-ORG>/qwen2.5-coder-7b-sft-bokeh-v1` (SFT)
- `<ANON-ORG>/qwen2.5-coder-7b-ga-bokeh-v1` (GA)
- `<ANON-ORG>/qwen2.5-coder-7b-relearn-bokeh-v1` (ReLearn)
- `<ANON-ORG>/qwen2.5-coder-7b-tv-fisher-bokeh-v1` (TV-Fisher; PDR-destroyed at evaluation)

*Qwen2.5-Coder-7B × Pydantic cell:*

- `<ANON-ORG>/qwen2.5-coder-7b-sft-pydantic-v1` (SFT)
- `<ANON-ORG>/qwen2.5-coder-7b-ga-pydantic-v1` (GA)
- `<ANON-ORG>/qwen2.5-coder-7b-relearn-pydantic-v1` (ReLearn)

*CodeLlama-7B × Pydantic cell:*

- `<ANON-ORG>/codellama-7b-sft-pydantic-v1` (SFT)
- `<ANON-ORG>/codellama-7b-ga-pydantic-v1` (GA)
- `<ANON-ORG>/codellama-7b-relearn-pydantic-v1` (ReLearn)

The Fisher-weighted task-vector merge is computed at evaluation time from the two TV adapters; no merged TV-Fisher checkpoint is shipped because the corrected  $4 \times 6$  sweep in §4 produced no working operating point.

**Data artifacts.** The v3 forget-set, retain-set, attack suite, and L5 probing dataset are distributed under data/v3/ in this repository: ga\_forget\_qa\_v3.json (94 entries), ga\_relearn\_retain\_v3.json (373 entries), relearn\_forget\_augmented\_v3.jsonl (564 entries = 94 originals × 3 vague-answer + 94 originals × 3 adversarial-answer augmentations), sft\_retain\_v3.json (373 entries), level5\_probing\_dataset.jsonl (187 entries: 94 deprecated, 93 current). The detector and v3 builder scripts live at data\_generation/tests/check\_bokeh\_deprecations.py and data\_generation/build\_v3\_unlearning\_data.py. Pydantic-specific data artifacts are distributed under data/v3/ with the \_pydantic suffix: sft\_retain\_pydantic\_v3.json (518 entries), ga\_forget\_qa\_pydantic\_v3.json (297 entries), relearn\_forget\_augmented\_pydantic\_v3.jsonl (1782 entries = 297 × 6 augmentations), and forget\_set\_pydantic.json (297 entries across 16 deprecation kinds).

**Seeds.** Training: seed = 3407 (Unsloth default) is fixed across SFT, GA, ReLearn, and the two TV adapter trainings. Fisher sweep single-pass: seed = 42. Bootstrap top-2 re-evaluation: seeds = {7, 13, 42} with do\_sample=True, temperature 0.7, top-p 0.95, per-seed random subsampling of evaluation prompts. L5 probe split: random\_state = 42, 80/20 stratified.

**Compute environment.** All training and evaluation runs were executed on Colab Pro A100 (40 GB). Software stack: PyTorch 2.x, transformers, peft, trl, Unsloth (Han et al., 2024), bitsandbytes for 4-bit (NF4) quantization (Detmers et al., 2023), scikit-learn for the L5 probe, huggingface\_hub for artifact upload, and transformers-bundled CodeBERT for similarity scoring. Approximate wall-clock per recipe: SFT 40 min, GA 25 min, ReLearn 60 min, each TV adapter 20 min, full Fisher rev3 sweep 7 min, Eval\_01 (5 models, 6 axes) 150 min, Eval\_01b (5 models, HumanEval + PDR) 40 min, Eval\_02 (5 models × 33 layers) 7 min. Cross-cell runs: approximately 300 min total per cell (6 training notebooks + 1 combined eval notebook).

**Live telemetry.** All long-running notebooks emit per-event JSONL telemetry to a shared Drive folder; the polling script scripts/poll\_telemetry.py renders progress without requiring stdout capture.

**Code release.** All project materials — source code for the six-level attack-suite generator, the v3 data builder, the deprecation detector catalog, the Fisher-sweep notebook (TV\_03 rev3), the representation-probe notebook (Eval\_02), and the figure-generation script — along with the five model artifacts listed above and the v3 data artifacts under data/v3/, are available at <https://anonymous.4open.science/r/codeunlearn-bench-67DD/> (anonymous mirror for review; de-anonymized repository and Hugging Face model namespace will be linked in the camera-ready version).

**Anti-patterns avoided.** Three implementation pitfalls cost us time during this study; we list them so others do not repeat them. (i) Gradient ascent on raw bnb-4bit base parameters is a silent no-op: the LoRA wrap must precede ascent, and p.grad.neg\_() must be applied explicitly post-backward (graph-level negation is swallowed by Unsloth’s checkpointed backward). (ii) Recovery fine-tuning at L4 must reactivate the existing adapter via FastLanguageModel.for\_training(model); calling get\_peft\_model stacks a second adapter and silently changes what is being measured. (iii) In-distribution validity guards calibrated on the forget-set distribution can pass spuriously on destroyed checkpoints; a held-out general-code generation spot-check is required (see §4).