

Beyond Forget Scores: Auditing Quantization-Robust LLM Unlearning with Q-ROU

Anonymous ACL submission

Abstract

LLM unlearning is often judged by whether a target answer disappears, but deployed models fail in richer ways: the same edit can collapse semantic neighbors, induce degenerate generation, reintroduce forgotten content after INT4/NF4 compression, or recover it after benign post-edit training. We therefore argue that post-deployment unlearning should be evaluated as an audit problem, and introduce Quantization-Robust Orthogonal Unlearning (Q-ROU), a compact operating point combining active semantic-neighbor retention with bounded KL-to-uniform forgetting plus SLUG/QuantNoise deployment stabilizers. On a 28-probe 3B multi-entity stress test, Q-ROU reaches 27/28 in FP16 and 28/28 in INT4, while non-AR baselines collapse neighbor retention. On a 65-probe expanded suite, Q-ROU achieves 25/25 target suppression with 17/20 neighbor retention in FP16; aggressive baselines either collapse neighbors or reach stronger token suppression only in degenerate generation regimes. Real NF4 checks preserve the main Q-ROU operating point, and LoRA-concentrated baselines show the key distinction: low-bit survival alone does not guarantee semantic-neighbor preservation. Standardized TOFU runs support the same failure-mode analysis, while MUSE news, same-work neighbors, and post-edit recurrence expose boundaries that target-only metrics hide. The resulting claim is limited to INT4/NF4 deployment audits: Q-ROU delivers selective, low-bit-stable suppression at a substantially stronger operating point than target-only metrics suggest.

1 Introduction

A high forget score can be the least informative part of a post-deployment unlearning result. Large language models (LLMs) memorize vast amounts of training data, including information that may need removal after deployment—for privacy regulations (e.g., GDPR “right to be forgotten” (European Par-

liament and Council of the European Union, 2016)), copyright concerns, or the emergence of harmful content. The natural first test is whether the target answer disappears, but that test is easy to pass for the wrong reasons. An edit can damage adjacent facts, make generation incoherent, fail after 4-bit deployment, or become easy to reverse through a small benign update. We use *deployment audit* to name the stronger requirement: post-hoc unlearning should be judged by the failure modes it rules out, not by a single target-side score.

Existing gradient-based methods (GA, GradDiff, REPBEND) share at least three deployment-critical limitations. First, they can incur substantial collateral damage to semantically adjacent knowledge: recent mechanistic analyses suggest that concept representations are distributed and can share nearby feature structure (Li et al., 2025), so updates aimed at a forget set can overwrite nearby capabilities. In our multi-entity stress test, the non-AR baselines GA, GradDiff, and RepBend collapse neighbor retention to 0/8. Second, they exhibit a token-generation tradeoff: aggressive objectives can suppress target tokens only by entering regimes of severe generation degeneration (Appendix D.4). Third, these methods can be brittle under deployment compression (Zhang et al., 2025; Abitante et al., 2026). Post-training INT4 quantization perturbs local decision margins and can cause partially forgotten information to resurface. Recent work (Abitante et al., 2026) demonstrates that standard full-parameter unlearning produces weight updates smaller than the quantization grid size, causing forgetting effects to be “masked” or “erased” during quantization—a phenomenon we term *forgetting regression*. To address these limitations, we propose **Quantization-Robust Orthogonal Unlearning (Q-ROU)**, a deployment-oriented operating point that combines bounded forgetting, explicit neighbor anchoring, restricted updates, and quantization-aware stabilization. *Active Retention*

085 (AR) uses a KL-divergence constraint to explic- 137
 086 itly preserve the model’s output distributions on 138
 087 neighbor-domain texts, directly targeting the se- 139
 088 mantic knowledge boundary rather than relying on 140
 089 generic retention data. *KL-to-uniform forgetting* 141
 090 bounds the forget objective by increasing target en- 142
 091 tropy, preventing optimization from entering degen- 143
 092 erate collapse regimes while still suppressing target 144
 093 facts. *Selective Layer Update Grouping (SLUG)* 145
 094 confines parameter modifications to knowledge- 146
 095 critical layers identified through activation analysis, 147
 096 reducing the surface area for collateral damage. Fi- 148
 097 nally, *QuantNoise injection* simulates quantization 149
 098 perturbations during training, steering the optimizer 150
 099 toward flat loss regions robust to post-deployment 151
 100 compression. Our primary contribution is a more 152
 101 demanding success criterion together with a method 153
 102 that reaches a strong deployment operating point 154
 103 under that criterion. Across the experiments be- 155
 104 low, methods that look successful on one axis rou- 156
 105 tinely fail on another: target-strong baselines dam- 157
 106 age neighbors, low-rank updates survive quantiza- 158
 107 tion without preserving semantic boundaries, and 159
 108 longer benchmark calibration can improve target 160
 109 suppression at unacceptable utility cost. Q-ROU 161
 110 is designed as a coherent operating point for this 162
 111 conjunction, with its scope made explicit by the 163
 112 audit suite rather than hidden behind a single forget 164
 113 score. While KL-based constraints for generic re- 165
 114 tention have been explored (e.g., SCRUB (Kurmanji 166
 115 et al., 2023), SalUn (Fan et al., 2024)), Q-ROU ex- 167
 116 plicitly anchors *semantic neighbors* (AR) and pairs 168
 117 this with a bounded KL-to-uniform forget objective 169
 118 to avoid degenerate regimes. Q-ROU combines 170
 119 a small core unlearning mechanism (AR + KL- 171
 120 to-uniform) with deployment-focused stabilizers 172
 121 (SLUG + QuantNoise) and lighter auxiliary regu- 173
 122 larizers (margin/orthogonality/EWC). This decom- 174
 123 position matters for interpreting the experiments: 175
 124 the paper does not claim that every term is equally 176
 125 essential, but that the layered combination reaches 177
 126 a stronger audited operating point than the com- 178
 127 pared alternatives. Concurrent work (Dorna et al., 179
 128 2025) likewise benchmarks multiple unlearning 180
 129 metrics, reinforcing our concern that keyword-style 181
 130 or membership-based summaries alone can miss 182
 131 failures exposed by adversarial or generation-level 183
 132 audits. We evaluate deployment unlearning along 184
 133 four audit axes—selective forgetting, deployment 185
 134 compression, extraction robustness, and post-edit 186
 135 persistence—because each axis exposes a differ-
 136 ent way for an apparent unlearning success to fail

in use. This positioning complements benchmark- 137
 standardization work such as TOFU (Maini et al., 138
 2024), MUSE (Shi et al., 2025), RWKU (Jin et al., 139
 2024), and OpenUnlearning (Dorna et al., 2025): 140
 our focus is to add deployment-specific failure 141
 modes that standard target/utility summaries can 142
 miss. It is motivated by recent benchmark-critique 143
 work showing that independent forget/retain evalua- 144
 tions can overstate progress and hide dependencies 145
 between what appears forgotten and what should 146
 remain intact (Thaker et al., 2025). It is also dis- 147
 tinct from recent low-rank quantization-robust un- 148
 learning (Abitante et al., 2026), reasoning-trace un- 149
 learning (Yoon et al., 2025), and relearning-attack 150
 defenses (Xiao et al., 2026): Q-ROU is not opti- 151
 mized for a single one of these axes, but audited for 152
 their conjunction under an explicit scope. It is also 153
 distinct from contemporaneous anonymous work 154
 on operator-level linear-access auditing in LLM 155
 representations (Anonymous, 2026a): that paper 156
 asks whether linear evidence supports readout, con- 157
 trol, and transfer claims in representation analy- 158
 sis, whereas Q-ROU targets post-deployment un- 159
 learning under suppression, neighbor preservation, 160
 low-bit stability, and bounded recurrence. Like- 161
 wise, contemporaneous anonymous work on pre- 162
 dicting collateral damage from Gradient Ascent 163
 unlearning (Anonymous, 2026b) studies which non- 164
 target concepts are most vulnerable under GA up- 165
 dates, framing the problem as a mechanistic ranking 166
 task; Q-ROU, by contrast, proposes a deployment- 167
 audited unlearning operating point and audit suite 168
 aimed at reducing such damage rather than predict- 169
 ing it. 170

2 Methodology 171

2.1 Problem Formulation 172

Let θ denote the parameters of a pre-trained LLM 173
 and $P_\theta(y|x)$ the conditional probability distribu- 174
 tion over output tokens y given input x . We define 175
 three evaluation sets: $\mathcal{D}_{\text{target}}$ (knowledge to be re- 176
 moved), $\mathcal{D}_{\text{neighbor}}$ (semantically related knowledge 177
 to be preserved), and $\mathcal{D}_{\text{general}}$ (domain-independent 178
 knowledge to be preserved). The goal of selec- 179
 tive unlearning is to find parameters θ^* such that, 180
 for each target context $x \in \mathcal{D}_{\text{target}}$ and each desig- 181
 nated target token position $t \in S_f(x)$ with target 182
 token y_t , the model suppresses $P_{\theta^*}(y_t | x, y_{<t})$, 183
 while maintaining $P_{\theta^*}(y|x) \approx P_\theta(y|x)$ for inputs 184
 in $\mathcal{D}_{\text{neighbor}} \cup \mathcal{D}_{\text{general}}$. Here $\mathcal{D}_{\text{target}}$ denotes target 185
 prompts/contexts, and forgetting is enforced only 186

Audit axis	Evidence	Primary readout	Scope limit
Selective forgetting	28-probe core, 65-probe expanded suite, and focused 3-seed confirmations in FP16/INT4	T/N/G pass counts with token-probability audits	strongest evidence is on controlled custom suites; same-work neighbors remain the hardest frontier
Deployment compression	Matched FP16, fake INT4, and real NF4 checks on the same suites	T/N/G retention plus Δ_{zombie}	establishes low-bit stability at the tested operating points, not a universal PTQ guarantee
Extraction robustness	Six-protocol / 55-probe depth suite plus budgeted extraction sweeps, with 160-step extension for extraction only	protocol pass rates and attack-budget leakage	bounded black-box extraction families only
Post-edit persistence	Bounded jog audits on Llama, Qwen, and Phi with supporting 3-seed checks	worst post-jog target count	bounded recurrence evidence on the tested jog families

Audit design overview. Four deployment failure modes, their evidence, primary readouts, and scope limits.

187 on designated target-answer token positions within
188 those contexts. Unlike data-partitioning approaches
189 such as SISA training (Bourtole et al., 2021),
190 which achieve unlearning by retraining isolated
191 model shards, Q-ROU directly modifies parameters
192 of a deployed model to enable post-hoc knowledge
193 removal without access to the original training data
194 or retraining infrastructure. Additionally, we re-
195 quire post-quantization robustness. For neighbor
196 and general sets, the output distributions should
197 remain stable under a post-training quantization op-
198 erator $Q(\cdot)$ (e.g., INT4), meaning $P_{Q(\theta^*)}(y|x) \approx$
199 $P_{\theta^*}(y|x)$. For the target set, the key requirement is
200 the absence of forgetting regression after quantiza-
201 tion. We formalize this target-side behavior through
202 the *zombie delta* metric: $\Delta_{\text{zombie}} = \mathbb{E}[\delta(x, y)]$,
203 where $\delta(x, y) = P_{Q(\theta^*)}(y | x) - P_{\theta^*}(y | x)$ and y
204 denotes the target token. Positive Δ_{zombie} indicates
205 *zombie knowledge* (forgotten information resurfac-
206 ing under quantization), while negative values indi-
207 cate that low-bit conversion does not cause forget-
208 ting regression and may further suppress borderline
209 targets.

210 2.2 Q-ROU Framework

211 Q-ROU optimizes a composite objective function
212 over K training steps, modifying only parameters
213 in a selected subset of layers (SLUG). We separate
214 the method into four layers of responsibility. **Core**
215 **unlearning** consists of KL-to-uniform forgetting
216 plus AR; **deployment stabilizers** consist of SLUG
217 and QuantNoise; **auxiliary regularizers** consist
218 of margin, orthogonality, and EWC; and **PTI** is a
219 post-edit persistence extension evaluated separately
220 rather than part of the base method. The objective
221 therefore follows a two-tier training design. The
222 **Core Q-ROU** objective handles the fundamental

unlearning task:

$$\mathcal{L}_{\text{core}} = \lambda_f \mathcal{L}_{\text{forget}} + \lambda_r \mathcal{L}_{\text{retain}} + \lambda_a \mathcal{L}_{\text{active}} \quad (1) \quad 224$$

225 To improve deployment robustness under low-bit
226 quantization, we complement this with an auxiliary
227 structural regularization suite:

$$\mathcal{L}_{\text{aux}} = \lambda_m \mathcal{L}_{\text{margin}} + \lambda_o \mathcal{L}_{\text{ortho}} + \mathcal{L}_{\text{EWC}} \quad (2) \quad 228$$

229 The total loss is $\mathcal{L} = \mathcal{L}_{\text{core}} + \mathcal{L}_{\text{aux}}$. As the ablation
230 results later show (Appendix C.6), the auxiliary
231 suite (\mathcal{L}_{aux}) should be read as robustness support
232 rather than as the conceptual core of the method: at
233 larger scales and easier settings it has little effect
234 on top-level pass counts, but it remains useful on
235 harder boundary cases and in low-bit audits.

Forget and Retain Objectives. The forget loss
236 uses a bounded KL-to-uniform objective that sup-
237 presses confident target-token predictions by in-
238 creasing output entropy on target contexts: $\mathcal{L}_{\text{forget}} =$
239 $\frac{1}{|S_f|} \sum_{t \in S_f} D_{\text{KL}}(P_{\theta}(\cdot | x, y_{<t}) \| U)$, where S_f
240 denotes the set of target token positions and U is the
241 uniform distribution over the vocabulary. Mini-
242 mizing this objective is equivalent, up to the con-
243 stant $\log |\mathcal{V}|$, to maximizing the output entropy
244 $H(P_{\theta})$ at the targeted positions. A passive re-
245 tain loss preserves general model behavior using a
246 KL-divergence constraint on general-domain texts:
247 $\mathcal{L}_{\text{retain}} = \mathbb{E}_{x_g \sim \mathcal{D}_{\text{general}}} [D_{\text{KL}}(P_{\theta_0}(\cdot | x_g) \| P_{\theta}(\cdot | x_g))]$,
248 where θ_0 denotes the frozen pre-unlearning param-
249 eters. 250

Active Retention (AR). Active Retention is the
251 core mechanism for preserving neighbor knowledge.
252 Unlike traditional methods that rely on passive re-
253 taining via generic corpora (which the literature
254 shows is insufficient against knowledge entangle-
255 ment), AR explicitly defends the semantic frontier 256

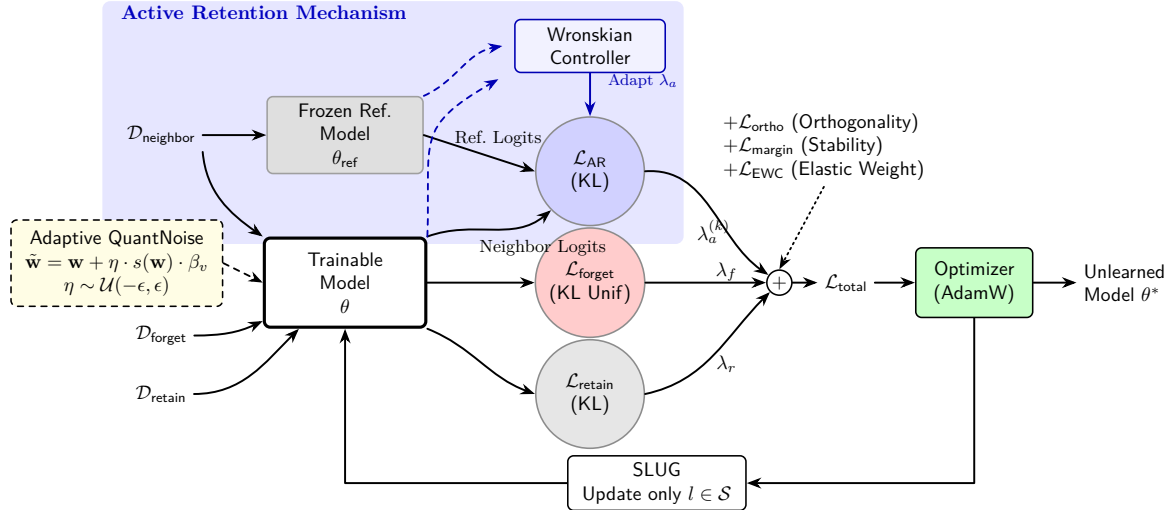


Figure 1: Q-ROU method overview. Target contexts optimize a bounded KL-to-uniform forget loss, semantic-neighbor and general retain sets are anchored to the frozen reference model through KL constraints, and the combined objective updates only the SLUG layer subset while QuantNoise and auxiliary regularizers support low-bit deployment stability. The Wronskian controller shown in the diagram is an optional diagnostic/adaptive variant; the main reported operating points use fixed coefficients unless otherwise stated.

257 around the forgotten domain by anchoring the out- 284
 258 put distribution on neighbor-domain texts to the 285
 259 frozen reference model: 286

$$260 \quad \mathcal{L}_{\text{active}} = \mathbb{E}_{x_n \sim \mathcal{D}_{\text{neighbor}}} \left[287 \right. \\
 261 \quad \left. D_{\text{KL}}(P_{\theta_0}(\cdot|x_n) \parallel P_{\theta}(\cdot|x_n)) \right] \quad (3) \quad 288$$

262 In its standard formulation, Q-ROU uses a fixed 291
 263 penalty coefficient λ_a for this objective. Supporting 292
 264 analysis in the Appendix shows that the KL penalty 293
 265 bounds neighbor drift through a Fisher-geometric 294
 266 constraint (Proposition 2). This Fisher-geometric 295
 267 interpretation shares conceptual foundations with 296
 268 prior work on linearized fine-tuning (Achille et al., 297
 269 2021), which uses local quadratic / Fisher-style ap- 298
 270 proximations to stabilize fine-tuning updates. We 299
 271 use fixed λ_a for all reported operating points; adap- 300
 272 tive variants are discussed only in the Appendix. 301

272 **Selective Layer Update Grouping (SLUG).** Mo- 302
 273 tivated by knowledge-localization and activation- 303
 274 pattern research (Meng et al., 2022, 2023; Dai et al., 304
 275 2022; Wang et al., 2025a), we restrict updates to 305
 276 knowledge-critical mid-to-deep layers, leaving shal-
 277 low layers frozen. This design shares conceptual
 278 foundations with weight saliency approaches (Fan
 279 et al., 2024), which selectively update important
 280 parameters to minimize collateral damage—though
 281 SLUG operates at the layer granularity rather than
 282 individual weights. In our experiments, we update
 283 5 layers per model; the selection procedure and

exact layer sets are reported in the Appendix. All 284
 parameters outside \mathcal{S} remain frozen during unlearn- 285
 ing, providing inherent protection for knowledge 286
 stored in non-SLUG layers and reduced compu- 287
 tational cost. This design shares conceptual simi- 288
 larities with LoRA-based unlearning (Abitante 289
 et al., 2026), which concentrates updates into low- 290
 rank subspaces to achieve quantization robustness. 291
 While LoRA achieves concentration through rank 292
 constraints on adapter matrices, SLUG achieves 293
 it through layer selection—both strategies amplify 294
 per-parameter update magnitudes to exceed quan- 295
 tization boundaries, though SLUG additionally lever- 296
 ages knowledge localization to minimize collateral 297
 damage. An Elastic Weight Consolidation (EWC) 298
 penalty (Kirkpatrick et al., 2017) further constrains 299
 these SLUG-layer parameters to remain close to 300
 their original values: 301

$$302 \quad \mathcal{L}_{\text{EWC}} = \frac{\lambda_{\text{ewc}}}{2} \sum_{l \in \mathcal{S}} \sum_i F_i (w_i - w_i^{(0)})^2 \quad (4) \quad 303$$

where $w_i^{(0)}$ denotes the original pre-unlearning 303
 weights and F_i is the Fisher information approx- 304
 imated from retain data. 305

QuantNoise and Regularization. During training, 306
 we inject uniform unstructured noise into train- 307
 able SLUG-layer weights: $\tilde{w} = w + \eta$, where 308
 $\eta \sim \mathcal{U}(-\epsilon_q, \epsilon_q)$. Intuitively, this is not “freez- 309
 ing” parameters; it actively flattens the local loss 310
 basin so INT4 rounding noise is absorbed instead 311

312 of undoing forgetting updates. We provide a formal
313 derivation in Appendix A showing that inde-
314 pendent coordinate noise translates to an implicit
315 weighted Hessian-trace penalty, $\text{tr}(H\Sigma_\eta)$, with the
316 ordinary trace recovered only in the isotropic-noise
317 special case. We further add a margin loss to main-
318 tain a safety buffer under quantization error, and
319 an orthogonality penalty motivated by evidence
320 that truth-related structure can occupy broad lin-
321 ear directions in representation space (Marks and
322 Tegmark, 2023). Additional details and optional
323 variants are reported in the Appendix.

324 Analytical Framework and Training Procedure.

325 Two lightweight propositions formalize our design:
326 (1) when the log-probability margin exceeds quanti-
327 zation error, pass/fail decisions are preserved under
328 INT4 rounding (Appendix Proposition 1); (2) AR is
329 the Lagrangian relaxation of constrained neighbor-
330 drift minimization, where the KL constraint gives
331 a local Fisher-geometric interpretation of neighbor
332 protection: under the usual small-update quadratic
333 KL approximation, $\Delta\theta^\top \mathcal{I}_{\text{NBR}} \Delta\theta$ is controlled by
334 the AR multiplier, biasing updates toward direc-
335 tions that leave neighbor outputs invariant (Ap-
336 pendix Proposition 2). Full statements, assump-
337 tions, proofs, and the complete training algorithm
338 (Algorithm 1) are in Appendix A.

339 3 Experimental Setup

340 **Models and Hyperparameters.** We primarily
341 evaluate on Llama-3.2-3B (Meta AI, 2024) for
342 the custom-probe audits. We additionally include
343 Llama-3.1-8B for TOFU standardized benchmarks
344 and Qwen2.5-3B / Phi-3-mini-4k-instruct for sup-
345 porting cross-architecture transfer and persistence
346 checks discussed below and in the appendix. For
347 the custom-probe experiments, we optimize Q-ROU
348 with AdamW for up to 40 steps, enabling efficient
349 post-deployment edits without prolonged retraining.
350 For standardized TOFU and long-horizon extrac-
351 tion analyses, we additionally evaluate 160-step and
352 longer schedules as specified in the corresponding
353 sections and appendix. Unless otherwise specified,
354 the primary loss coefficients on Llama-3.2-3B are
355 set to $\lambda_f = 2$, $\lambda_r = 15$, $\lambda_a = 80$, $\lambda_m = 5$, $\lambda_o = 5$,
356 and $\lambda_{\text{ewc}} = 1$. For the final calibrated 40-step multi-
357 entity operating point and the supporting LoRA
358 / memory-jog audits, we increase only the forget
359 coefficient to $\lambda_f = 20$ after a calibration sweep
360 showed that $\lambda_f = 2$ under-suppresses the 13 tar-
361 get probes while leaving the retention coefficients

unchanged (Appendix C.3).

362

Datasets and Evaluation Protocol. The primary
363 stress test is multi-entity forgetting across two well-
364 known fictional domains with distinct neighbor and
365 general retain sets. We define $\mathcal{D}_{\text{target}}$ as 18 derived
366 target contexts/prompts spanning Harry Potter and
367 LOTR; $\mathcal{D}_{\text{neighbor}}$ as 11 derived contexts/prompts
368 from Star Wars; and $\mathcal{D}_{\text{general}}$ as 20 non-fiction de-
369 rived contexts/prompts. Evaluation uses binary
370 pass counts across three categories: target (TGT,
371 13 probes), neighbor (NBR, 8), and general (GEN,
372 7). A target probe passes if its predicted probability
373 $P < 0.01$; neighbor/general pass if the probability
374 ratio vs. baseline ≥ 0.30 (fallback: $P > 0.05$). We
375 prioritize binary pass-count metrics over continu-
376 ous averages to prevent a single severe knowledge
377 leak from being masked by widespread but shallow
378 forgetting. The strict $P < 0.01$ margin provides a
379 conservative buffer against low-bit drift, aiming to
380 keep target probabilities below $P < 0.05$ after preci-
381 sion degradation (Appendix B.3). Because thresh-
382 olded counts can be brittle on small probe sets, we
383 pair them with raw probabilities, zombie deltas,
384 threshold-sensitivity sweeps, and Wilson/bootstrap
385 intervals in the Appendix; no claim in the paper
386 relies on a single thresholded table alone. Beyond
387 the 28-probe core set, we evaluate on a 65-probe
388 expanded set (paraphrased and indirect references),
389 six adversarial depth protocols, and budgeted ex-
390 traction.

391

Baselines and Precision Settings. We com-
392 pare Q-ROU against Gradient Ascent (GA (Jang
393 et al., 2023; Yao et al., 2024)), GradDiff
394 (GRADDIFF (Maini et al., 2024)), and Repre-
395 sentation Bending (REPBEND (Yousefpour et al.,
396 2025)). We also include comparisons against
397 NPO (Zhang et al., 2024), RMU (Li et al., 2024),
398 and GRU (Wang et al., 2025b) in our expanded eval-
399 uation. We also evaluate AR-augmented variants
400 (*e.g.*, GA+AR) for mechanism isolation. To test
401 the closest low-rank quantization-robust competi-
402 tor, we add LoRA-concentrated baselines: LoRA-
403 GA, LoRA-GA+AR, and LoRA-KL+AR, using
404 rank-16 adapters over attention and MLP projec-
405 tions plus a small slug_down adapter control (Ap-
406 pendix C.3). All models are evaluated in FP16 pre-
407 cision and subsequently after post-training INT4
408 quantization (group size 32) to formally measure
409 post-quantization performance drift.

410

4 Results

4.1 Standardized Benchmark: TOFU

To evaluate standardized unlearning, we run the TOFU benchmark (160 steps, robust profile) on both Llama-3.2-3B and Llama-3.1-8B, and then extend the hardest 8B forget10 case with long-horizon sweeps in the appendix.

Table 1: TOFU results under the standard 160-step protocol (single seed). Long-horizon 8B forget10 sweeps and tuned multi-seed hard-split confirmations are reported in Appendix C.4.

Subset	Model	Method	Forget@0.01	Forget@0.05	Gen Leak	Utility	Truth
forget01	Llama-3.2-3B	Q-ROU	100.0%	100.0%	0.0%	98.0%	100.0%
	Llama-3.2-3B	GA+AR	50.0%	50.0%	5.0%	99.0%	100.0%
	Llama-3.1-8B	Q-ROU	100.0%	100.0%	0.0%	100.0%	95.83%
	Llama-3.1-8B	GA+AR	50.0%	50.0%	7.5%	97.0%	95.83%
forget05	Llama-3.1-8B	Q-ROU	100.0%	100.0%	0.0%	69.0%	95.83%
	Llama-3.1-8B	GA+AR	11.0%	12.0%	34.5%	98.0%	100.0%
forget10	Llama-3.1-8B	Q-ROU	72.0%	80.25%	7.0%	82.0%	100.0%
	Llama-3.1-8B	GA+AR	5.75%	7.0%	36.75%	100.0%	91.67%

Q-ROU outperforms GA+AR on the forget/leakage metrics across all subsets. On forget01 and forget05, it achieves full 100% Forget@0.01 with 0% leakage. On forget10, the fixed-budget evaluation maintains a wide gap (72.0% vs. 5.75% Forget@0.01; 7.0% vs. 36.75% leakage). We report GA+AR in the main text as the strongest selective baseline under our retention-oriented criteria; additional baseline results are provided in the Appendix. At the same fixed 160-step budget, 8B utility on forget05/forget10 drops relative to GA+AR, reflecting a forgetting–utility trade-off at harder targets; we analyze longer-horizon behavior and ablation sensitivity in the Appendix. Longer-horizon 8B sweeps (Appendix C.4) reveal the 160-step forget10 checkpoint is materially under-saturated: Q-ROU enters a high-performing regime by steps 288–416. End-to-end MUSE/WMDP runs complete without instability and clarify where longer-horizon calibration matters most. WMDP and MUSE books act as stable high-utility checks. MUSE news is the hardest external setting: the default 40-step row is under-saturated, a 120-step short-sequence calibration sweep reaches positive KnowMem reduction at 0.9934 utility, and a full-sequence 120/160/240-step continuation preserves high utility while showing that the benchmark’s full-sequence news forget metric remains difficult. We therefore use MUSE news as a boundary and calibration diagnostic, while the custom multi-entity audits remain the cleanest place to isolate Q-ROU’s selective-forgetting mechanism (Appendix B.5).

4.2 Audit Axis 1: Selective Forgetting

Beyond standardized benchmarks, we conduct the main diagnosis using a custom multi-entity suite designed to stress selective forgetting rather than target suppression alone. Table 2 reports the core operating points. In the multi-entity setting, Q-ROU reaches 27/28 in FP16 and 28/28 in INT4, outperforming the strongest baseline family under the same evaluation protocol. Notably, INT4 performance matches or slightly exceeds FP16 on this operating point, indicating that low-bit conversion does not induce forgetting regression here.

Table 2: Core pass-count results on Llama-3.2-3B (28-probe core set).

Setting	FP16 (T/N/G)	INT4 (T/N/G)
Single (Q-ROU)	23/23 (7/7, 9/9, 7/7)	23/23 (7/7, 9/9, 7/7)
Multi (Q-ROU)	27/28 (13/13, 7/8, 7/7)	28/28 (13/13, 8/8, 7/7)
Multi (GA+AR)	22/28 (7/13, 8/8, 7/7)	22/28 (7/13, 8/8, 7/7)

Baseline GA/GRADDIFF/REP BEND runs suffer *near-complete collapse* of neighbor retention (falling to 0/8) and damage general behavior. The key distinction is not target suppression alone, but whether suppression remains selective and stable under deployment precision changes. This same point appears in the matched LoRA audit: broad LoRA-GA survives low-bit conversion but still collapses semantic neighbors, whereas AR/KL LoRA variants preserve neighbors only by under-forgetting. We therefore treat low-rank concentration as supporting evidence for the compression story, not as a substitute for neighbor-aware suppression (Appendix C.3).

To test broader generalization, we scale up to a 65-probe expanded set (paraphrased and indirect references) and compare against seven baselines (Table 3). Appendix Table 43 then reproduces the focused five-method selective subset used for the detailed 40-step operating-point discussion. Here we report only the converged 40-step operating points used for the main claim. Appendix Table 27 gives the 20/40/60 step sweep and shows the step dependence directly: Q-ROU is still under-suppressed at 20 steps, cleanly overtakes GA+AR by 40 steps, and remains ahead at 60 steps even as GA+AR partially closes the target gap. Appendix Table 44 then confirms that this 40-step ordering is not a single-seed artifact: across three seeds, Q-ROU stays at 25/25 target suppression in both FP16 and INT4, while GA+AR remains fixed at 13/25 target pass. Appendix D.4 further reports learning-rate, step-count,

494 and AR-weight controls for the strongest selective
 495 baseline family. Extending aggressive baselines
 496 without additional stabilizers often degrades gener-
 497 ation, which we treat as a failure mode rather than
 498 a valid operating point. Q-ROU is the only method
 499 achieving 25/25 target removal in FP16 while still
 500 retaining 17/20 neighbor pass counts, and it remains
 501 the only method with 25/25 target removal together
 502 with $\geq 15/20$ neighbor retention under INT4. All
 503 other high-target methods (GA, RepBend, RMU,
 504 GRU) achieve $\leq 1/20$ neighbor retention. NPO pre-
 505 serves neighbors (15/20) but achieves only 15/25
 506 target removal. Because aggressive GA-family set-
 507 tings can trade token-level suppression for genera-
 508 tion degeneration, we interpret this pass-count table
 509 jointly with the operating-point audit in Table 5 and
 510 the generation-level analyses in Appendix D.4. Ad-
 511 ditional plug-in baseline experiments are reported
 512 in the Appendix.

Table 3: Seven-method baseline comparison on the 65-probe expanded set (Llama-3.2-3B). Q-ROU and GA+AR are reported at their 40-step coherent operating points; other baselines follow their published recommended budgets (20 steps). The 20/40/60 sweep is reported separately in Appendix Table 27. Extending aggressive baselines beyond 20 steps without stabilizers causes generation degeneration (Appendix D.4), which we treat as a failure mode. Q-ROU uniquely combines complete target removal with robust retention across FP16 and INT4; aggressive baselines (GA, RepBend, RMU, GRU) collapse neighbor knowledge. FP and I4 denote FP16 and INT4 precision respectively.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z
Original model	5/25	20/20	20/20	5/25	18/20	20/20	-0.055
Q-ROU (40 steps)	25/25	17/20	20/20	25/25	16/20	18/20	-0.018
GA+AR	13/25	18/20	20/20	13/25	15/20	19/20	-0.026
GA	25/25	0/20	0/20	25/25	0/20	0/20	+0.000
RepBend	25/25	0/20	12/20	25/25	0/20	11/20	-0.013
RMU	25/25	0/20	13/20	25/25	0/20	13/20	-0.002
NPO	15/25	15/20	18/20	16/25	15/20	17/20	-0.027
GRU	25/25	0/20	1/20	25/25	0/20	1/20	+0.000

513 Three-seed confirmation on the 28-probe suite
 514 and threshold sweeps on the larger audit support the
 515 same reading: the joint T/N/G advantage is stable
 516 across seeds, and the target-side Q-ROU conclu-
 517 sion remains stable around the default cutoff while
 518 GA+AR shifts more sharply with τ_f (Appendix
 519 Tables 44 and 50).

520 The same selective-forgetting story becomes
 521 sharper under a stricter neighbor definition. Split-
 522 ting the original neighbor set into alias, same-work,
 523 inter-domain, and general strata shows that same-
 524 work retention is the hardest frontier, not a con-
 525 tradiction of the main result. On this harder au-

dit, the default Q-ROU point still deletes the tar-
 get completely, while a fine-grained AR operat-
 ing point raises same-work retention from 9/36 to
 19/36 in FP16 and to 18/36 under three-seed BnB
 NF4 (Appendix Tables 11 and 58). The larger 81-
 probe copyright-like benchmark shows the same
 pattern: Q-ROU+FG-AR lifts same-work retention
 to 14.7/24 in FP16 and 12.7/24 under real NF4
 while staying target-strong (Appendix Table 12).

4.3 Audit Axis 2: Deployment Compression and Baseline Fairness

Ablation of Adaptive Mechanisms. Four-way
 ablation of the adaptive components (Standard,
 Adaptive QN Only, Adaptive AR Only, Full Adap-
 tive) reveals only small pass-count differences on
 the 3B expanded set (Appendix Table 36). This is
 expected: the adaptive controller is designed to mit-
 igate *tail* deployment failures (e.g., threshold britt-
 leness and instability under precision changes) rather
 than to monotonically improve mean pass counts,
 so its value is better interpreted alongside robust-
 ness audits (Appendix C.12). Notably, Q-ROU’s
 multi-component design addresses the 4-bit quan-
 tization forgetting regression challenge identified
 by (Abitante et al., 2026): while their LoRA-based
 recipe preserves forgetting through low-rank up-
 date concentration, Q-ROU combines SLUG (layer-
 wise concentration), QuantNoise (loss-landscape
 flattening), and structural regularizers to stabilize
 the selective-forgetting operating point beyond up-
 date magnitude alone.

Active Retention Transplants and Feature Removals.

The transplant results (Table 4) are an
early-phase mechanism diagnostic: we inspect a
 deliberately under-saturated 20-step checkpoint to
 separate AR’s boundary-anchoring role from even-
 tual target saturation. At that matched early point,
 removing AR from Q-ROU already collapses neigh-
 bor retention from 8/8 to 2/8 (FP16) and 0/8 (INT4)
 while accelerating target removal. Conversely,
 transplanting AR into aggressive baselines prevents
 neighbor collapse (0/8 \rightarrow 8/8), but this *neighbor-*
only sufficiency does not resolve the broader de-
 ployment failure modes (token-generation tradeoff
 and low-bit stability) that separate Q-ROU from
 GA+AR at the coherent 40-step operating points
 (Table 3; Appendix D.4).

Removing QuantNoise produces negligible FP16
 changes but reveals why low-bit robustness matters:
 FP TGT shifts from 7/13 to I4 TGT 9/13 (+2) while

576 NBR stays at 8/8. In other words, the target bound-
577 ary is still moving under quantization even when
578 the headline count improves. AR bounds collateral
579 damage, while SLUG and QuantNoise make the
580 operating point less sensitive to that low-bit pertur-
581 bation. EWC and the remaining regularizers are
582 better read as auxiliary support than as the concep-
583 tual core.

Table 4: Early-phase AR mechanism isolation and ablation (Llama-3.2-3B, multi-entity, 20 steps). This matched checkpoint separates mechanism roles before convergence. The full 40-step operating-point separation is reported in Table 3 and Appendix Table 18.

Configuration	FP16 (T/N/G)	INT4 (T/N/G)
GA	13/13, 0/8, 0/7	13/13, 0/8, 0/7
GA+AR	4/13, 8/8, 7/7	6/13, 8/8, 7/7
Q-ROU	4/13, 8/8, 7/7	4/13, 8/8, 7/7
Q-ROU -AR	13/13, 2/8, 7/7	13/13, 0/8, 7/7
Q-ROU -QuantNoise	7/13, 8/8, 7/7	9/13, 8/8, 7/7

584 Scale-Dependent Component Contribution.

585 Component ablation at 3B scale reveals a different
586 sensitivity landscape from 0.5B. Removing
587 Orthogonality, EWC, or Margin individually yields
588 *negligible* pass-count changes ($|\Delta\bar{p}_{TGT}| \leq 0.013$),
589 contrasting with 0.5B where the same removals
590 cause visible degradation (Appendix C.6). The
591 correct reading is that these terms are auxiliary
592 guards whose effect is most visible on harder
593 or more deployment-sensitive slices, including
594 loss-landscape sharpness and threshold sensitivity
595 (Appendix D.16; Appendix C.12).

596 The main takeaway is visible in Tables 3 and 6:
597 Q-ROU’s target suppression does not regress after
598 fake INT4 or real NF4 conversion, whereas target-
599 strong baselines achieve this only by sacrificing
600 retention or generation quality.

601 This resolves the main baseline-fairness ambi-
602 guity: GA+AR can be tuned toward stronger to-
603 ken suppression, but the matched 13/13 row is not
604 deployment-valid because its apparent “no leak”
605 generation is achieved by degeneration rather than
606 by coherent non-target behavior (Appendix D.4).
607 Accordingly, the GA+AR rows used elsewhere in
608 the main text are the strongest non-degenerate op-
609 erating points from this audit rather than the best
610 token-only checkpoints.

611 Real NF4 supports the same reading as the INT4
612 audit: the deployment ordering survives low-bit
613 conversion, even though small same-work drifts
614 remain on the hardest slices.

Table 5: Baseline-fairness audit on the 28-probe core suite (Llama-3.2-3B, FP16, 40 steps). High-LR GA+AR reaches 13/13 only via repetitive degeneration.

Method	LR / λ_a	TGT	NBR	GEN	Deg.	Reading
Q-ROU	main	13/13	7/8	7/7	0/7	coherent
GA+AR	$10^{-4}/160$	9/13	8/8	7/7	0/7	coherent; still leaking
GA+AR	$5 \times 10^{-4}/40$	13/13	8/8	7/7	7/7	repetitive collapse

Table 6: Main real-quantization check on the Llama-3.2-3B multi-entity suite. Q-ROU preserves the operating-point ordering under fake INT4 and deployed NF4, whereas GA+AR remains under-suppressed.

Method	FP16	Fake INT4	NF4
Q-ROU 40s	13/13, 7/8, 7/7	13/13, 8/8, 7/7	13/13, 7/8, 7/7
GA+AR 40s	7/13, 8/8, 7/7	7/13, 8/8, 7/7	7/13, 8/8, 7/7

The matched low-rank audit shows the same sep- 615
616 aration: broad LoRA-GA survives low-bit conver- 617
618 sion but still collapses neighbors, while LoRA- 619
620 GA+AR and LoRA-KL+AR preserve neighbors 621
622 only by under-forgetting (Appendix C.3). 623

624 4.4 Audit Axis 3: Extraction Robustness 625

626 Across six depth protocols, Q-ROU achieves 627
628 55/55 suppression in both FP16 and INT4, while 629
630 GA+AR leaves exploitable vulnerabilities. Chain- 631
632 of-Thought and budgeted extraction show the same 633
634 pattern: the 40-step point still leaves a 33.3% 635
636 INT4 budget channel, but the 160-step extension 637
638 closes all tested budgets to 0.0% (Appendix Ta- 639
640 bles 45, 52, 53). 641

642 4.5 Audit Axis 4: Post-Edit Persistence 643

644 Under bounded non-target memory-jog updates, 645
646 PTI acts as a useful Llama-family hardening step: 647
648 the worst post-jog target count improves from 7/13 649
650 to 12/13 in the narrow audit, from 6/13 to 8/13 in 651
652 the broader five-mode audit, and from 5/13 to 10/13 653
654 under matching deployed-NF4 evaluation. Qwen 655
656 and Phi remain feasible after architecture-aware re- 657
658 tuning but do not show the same count-level PTI 659
660 gain (Appendix C.3). 661

662 5 Conclusion 663

664 Q-ROU is best read as deployment-audited suppres- 665
666 sion rather than a one-number forget score. Across 667
668 the Llama-family audits, it suppresses targets while 669
670 preserving semantic neighbors and avoiding low-bit 671
672 regression. It reaches a stronger joint operating 673
674 point than the tested baselines. 675

646 **Limitations**

647 Our strongest evidence comes from controlled
648 custom-probe audits on sub-3B Llama-family mod-
649 els, with standardized TOFU runs extending to 8B.
650 The 28- and 65-probe suites are designed to isolate
651 failure modes cleanly, and the benchmark results
652 show that operating-point sensitivity remains real,
653 especially on MUSE news. The hardest remain-
654 ing selectivity issue is same-work retention. Fine-
655 grained audits show clear gains from stronger neigh-
656 bor anchoring, but even the more protective AR op-
657 erating point does not fully preserve all same-work
658 non-target facts. Hyperparameters also remain
659 architecture-dependent. Architecture-aware trans-
660 fer to Qwen and Phi is feasible, but the strongest PTI
661 recurrence gains appear on Llama-family models.
662 Benign post-deployment updates can still reactiv-
663 ate a subset of target probes after base unlearning,
664 which is why we report PTI as a useful hardening ex-
665 tension rather than fold it into the base method. Our
666 low-bit robustness evidence covers fake INT4 and
667 real NF4; broader GPTQ/AWQ-style post-training
668 quantizers remain an important next check. Finally,
669 our adversarial CoT and budgeted extraction probes
670 are bounded elicitation/extraction audits with oracle
671 risk-ranked variants, not a full reasoning-model un-
672 learning benchmark. R-TOFU-style step-wise trace
673 evaluation is a natural next extension, especially for
674 models whose safety-critical behavior depends on
675 hidden or explicit reasoning traces.

676 **Ethical Considerations**

677 This work targets safer post-deployment models
678 by enabling selective removal of harmful, sensi-
679 tive, or copyrighted knowledge without requiring
680 expensive full retraining. However, we acknowl-
681 edge the dual-use nature of unlearning frameworks:
682 the same mechanisms could theoretically be mis-
683 used to conceal model provenance, selectively sup-
684 press legitimate information, or deliberately intro-
685 duce knowledge blind spots. We mitigate these
686 risks by emphasizing auditable evaluation proto-
687 cols, explicit reporting of retention trade-offs, and
688 transparent documentation of unlearning scope and
689 failure modes. All experiments rely on fictional-
690 domain evaluation sets derived from widely known
691 franchises, benchmark-style settings, and a supple-
692 mentary audit of publicly documented biographical
693 facts about public figures only. We do not redis-
694 tribute copyrighted source text. No private per-
695 sonal information is collected or released, and the

public-biographical audit is limited to public, non- 696
sensitive facts rather than a redistributable personal- 697
information dataset. 698

To support independent evaluation, the paper and 699
appendix report the model variants, optimization 700
settings, probe construction rules, thresholding cri- 701
teria, quantization settings, and hardware condi- 702
tions needed to interpret every table and figure. The 703
manuscript is intended to stand on its own: each re- 704
ported claim is tied to a directly described protocol 705
rather than to an external code release. AI assis- 706
tants were used only in a limited supporting role 707
for English expression polishing, routine coding as- 708
sistance, lightweight scripting, figure drafting, and 709
internal manuscript diagnostics. All substantive 710
research decisions, theoretical claims, experiment 711
design, result interpretation, and final manuscript 712
wording were reviewed and approved by the authors, 713
who take full responsibility for the paper. 714

References 715

- João Vitor Boer Abitante, Joana Meneguzzo Pasquali, 716
Luan Fonseca Garcia, Ewerton de Oliveira, Thomas 717
da Silva Paula, Rodrigo C. Barros, and Lucas S. 718
Kupssinski. 2026. [Quantization-robust LLM un- 719
learning via low-rank adaptation](#). *arXiv preprint*. 720
- Alessandro Achille, Aditya Golatkar, Avinash Ravichan- 721
dran, Marzia Polito, and Stefano Soatto. 2021. [Lqf: 722
Linear quadratic fine-tuning](#). In *Proceedings of the 723
IEEE/CVF Conference on Computer Vision and Pat- 724
tern Recognition (CVPR)*, pages 15729–15739. 725
- Anonymous. 2026a. Linear access in llm representa- 726
tions: An audit of readout, control, and transfer. Un- 727
der review. 728
- Anonymous. 2026b. Which neighbors suffer? predict- 729
ing collateral damage from gradient ascent in llm 730
unlearning. Under review. 731
- Lucas Bourtole, Varun Chandrasekaran, Christopher A. 732
Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu 733
Zhang, David Lie, and Nicolas Papernot. 2021. [Ma- 734
chine unlearning](#). In *2021 IEEE Symposium on Secu- 735
rity and Privacy (SP)*, pages 141–159. 736
- Thomas M. Cover and Joy A. Thomas. 2006. [Ele- 737
ments of Information Theory](#), 2nd edition. Wiley- 738
Interscience, Hoboken, NJ. 739
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao 740
Chang, and Furu Wei. 2022. [Knowledge neurons 741
in pretrained transformers](#). In *Proceedings of the 742
60th Annual Meeting of the Association for Compu- 743
tational Linguistics (Volume 1: Long Papers)*, pages 744
8493–8502. Association for Computational Linguis- 745
tics. 746

747	Chandler Davis and W. M. Kahan. 1970. The rotation of eigenvectors by a perturbation. iii . <i>SIAM Journal on Numerical Analysis</i> , 7(1):1–46.	804
748		805
749		806
750	Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. LLM.int8(): 8-bit matrix multiplication for transformers at scale . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 30318–30332.	807
751		808
752		809
753		810
754		811
755	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. QLoRA: Efficient fine-tuning of quantized LLMs . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 10088–10115.	812
756		813
757		814
758		815
759		816
760	Vineeth Dorna, Anmol Mekala, Wenlong Zhao, Andrew McCallum, J. Zico Kolter, Zachary C. Lipton, and Pratyush Maini. 2025. Openunlearning: Accelerating llm unlearning via unified benchmarking of methods and metrics . In <i>Advances in Neural Information Processing Systems 38 (Datasets and Benchmarks Track)</i> .	817
761		818
762		819
763		820
764		821
765		822
766	Cynthia Dwork and Aaron Roth. 2014. The algorithmic foundations of differential privacy . <i>Foundations and Trends® in Theoretical Computer Science</i> , 9(3–4):211–487.	823
767		824
768		825
769		826
770	European Parliament and Council of the European Union. 2016. Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data (General Data Protection Regulation) . Official Journal of the European Union. OJ L 119, 4.5.2016, pp. 1–88.	827
771		828
772		829
773		830
774		831
775		832
776		833
777		834
778	Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. 2024. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation . In <i>The Twelfth International Conference on Learning Representations (ICLR 2024)</i> , pages 53643–53673.	835
779		836
780		837
781		838
782		839
783		840
784	Jiahui Geng, Qing Li, Herbert Woisetschlaeger, Zongxiong Chen, Fengyu Cai, Yuxia Wang, Preslav Nakov, Hans-Arno Jacobsen, and Fakhri Karray. 2025. A comprehensive survey of machine unlearning techniques for large language models . <i>arXiv preprint</i> .	841
785		842
786		843
787		844
788		845
789	Antonio Ginart, Melody Guan, Gregory Valiant, and James Y. Zou. 2019. Making AI forget you: Data deletion in machine learning . In <i>Advances in Neural Information Processing Systems</i> 32.	846
790		847
791		848
792		849
793	Aditya Golatkar, Alessandro Achille, and Stefano Soatto. 2020. Eternal sunshine of the spotless net: Selective forgetting in deep networks . In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9301–9309.	850
794		851
795		852
796		853
797		854
798	Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models . In <i>Advances in Neural Information Processing Systems</i> 36, pages 17643–17668.	855
799		856
800		857
801		858
802		859
803		860
		861
	Yihuai Hong, Lei Yu, Haiqin Yang, Shauli Ravfogel, and Mor Geva. 2025. Intrinsic test of unlearning using parametric knowledge traces . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 19513–19535. Association for Computational Linguistics.	804
		805
		806
		807
		808
		809
	Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. 2023. Knowledge unlearning for mitigating privacy risks in language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14389–14408. Association for Computational Linguistics.	810
		811
		812
		813
		814
		815
		816
		817
	Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen, Kang Liu, and Jun Zhao. 2024. Rwku: Benchmarking real-world knowledge unlearning for large language models . In <i>Advances in Neural Information Processing Systems 37 (Datasets and Benchmarks Track)</i> , pages 98213–98263.	818
		819
		820
		821
		822
		823
		824
	Dhruva Karkada, Daniel J. Korchinski, Andres Nava, Matthieu Wyart, and Yasaman Bahri. 2026. Symmetry in language statistics shapes the geometry of model representations . <i>arXiv preprint</i> .	825
		826
		827
		828
	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2017. Overcoming catastrophic forgetting in neural networks . <i>Proceedings of the National Academy of Sciences</i> , 114(13):3521–3526.	829
		830
		831
		832
		833
		834
		835
		836
	Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions . In <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 1885–1894. PMLR.	837
		838
		839
		840
		841
	Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. 2023. Towards unbounded machine unlearning . In <i>Advances in Neural Information Processing Systems</i> 36, pages 1957–1987.	842
		843
		844
		845
	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, and 27 others. 2024. The WMDP benchmark: Measuring and reducing malicious use with unlearning . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 28525–28550. PMLR.	846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
	Yuxiao Li, Eric J. Michaud, David D. Baek, Joshua Engels, Xiaoqing Sun, and Max Tegmark. 2025. The geometry of concepts: Sparse autoencoder feature structure . <i>Entropy</i> , 27(4):344.	858
		859
		860
		861

862	Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper,	Akira Sakai and Yuma Ichikawa. 2026. Sign lock-in:	917
863	Nathalie Baracaldo, Peter Hase, Yuguang Yao,	Randomly initialized weight signs persist and bottle-	918
864	Chris Yuhao Liu, Xiaojun Xu, Hang Li, Kush R.	neck sub-bit model compression.	919
865	Varshney, Mohit Bansal, Sanmi Koyejo, and Yang		
866	Liu. 2025. Rethinking machine unlearning for large	Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Mal-	920
867	language models.	ladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke	921
868	<i>Nature Machine Intelligence</i> ,	Zettlemoyer, Noah A. Smith, and Chiyuan Zhang.	922
	7:181–194.	2025. Muse: Machine unlearning six-way evalua-	923
869	Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen	tion for language models.	924
870	Casper, and Dylan Hadfield-Menell. 2024. Eight	<i>In The Thirteenth International Conference on Learning Representations</i>	925
871	methods to evaluate robust unlearning in llms.	<i>(ICLR 2025)</i> , pages 27797–27818.	926
872	<i>arXiv preprint.</i>		
873	Pratyush Maini, Zhili Feng, Avi Schwarzschild,	Iliia Shumailov, Jamie Hayes, Eleni Triantafyllou,	927
874	Zachary Chase Lipton, and J. Zico Kolter. 2024.	Guillermo Ortiz-Jimenez, Nicolas Papernot, Matthew	928
875	TOFU: A task of fictitious unlearning for LLMs.	Jagielski, Itay Yona, Heidi Howard, and Eugene Bag-	929
876	<i>In Proceedings of the First Conference on Language</i>	dasaryan. 2024. Ununlearning: Unlearning is not	930
877	<i>Modeling (COLM 2024).</i>	sufficient for content regulation in advanced genera-	931
		tive ai.	932
		<i>arXiv preprint.</i>	
878	Samuel Marks and Max Tegmark. 2023. The geometry	G. W. Stewart and Ji-guang Sun. 1990. <i>Matrix Pertur-</i>	933
879	of truth: Emergent linear structure in large language	<i>bation Theory.</i> Academic Press, Boston.	934
880	model representations of true/false datasets.		
881	<i>arXiv preprint.</i>	Pratiksha Thaker, Shengyuan Hu, Neil Kale, Yash Mau-	935
		rya, Zhiwei Steven Wu, and Virginia Smith. 2025.	936
882	Kevin Meng, David Bau, Alex Andonian, and Yonatan	Position: LLM unlearning benchmarks are weak mea-	937
883	Belinkov. 2022. Locating and editing factual asso-	sures of progress.	938
884	ciations in gpt.	<i>In 2025 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML),</i>	939
885	<i>Advances in Neural Information Processing Systems</i> ,	pages 520–533.	940
	volume 35, pages 17359–17372.		
886	Kevin Meng, Arnab Sen Sharma, Alex Andonian,	Anvith Thudi, Gabriel Deza, Varun Chandrasekaran,	941
887	Yonatan Belinkov, and David Bau. 2023. Mass-	and Nicolas Papernot. 2022. Unrolling sgd: Under-	942
888	editing memory in a transformer.	standing factors influencing machine unlearning.	943
889	<i>In The Eleventh International Conference on Learning Representa-</i>	<i>In 2022 IEEE 7th European Symposium on Security and</i>	944
890	<i>tions (ICLR 2023).</i>	<i>Privacy (EuroS&P)</i> , pages 303–319. IEEE.	945
891	Meta AI. 2024. Llama 3.2: Revolutionizing edge ai	Haoyu Wang, Zhuo Huang, Xiaolong Wang, Bo Han,	946
892	and vision with open, customizable models.	Zhiwei Lin, and Tongliang Liu. 2026. Megu:	947
893	Meta AI Blog.	Machine-guided unlearning with target feature disen-	948
		tanglement.	949
		<i>arXiv preprint.</i>	
894	Markus Nagel, Marios Fournarakis, Rana Ali Amjad,	Yudong Wang, Damai Dai, Zhe Yang, Jingyuan Ma, and	950
895	Yelysei Bondarenko, Mart van Baalen, and Tijmen	Zhifang Sui. 2025a. Exploring activation patterns of	951
896	Blankevoort. 2021. A white paper on neural network	parameters in language models.	952
897	quantization.	<i>In Proceedings of the AAAI Conference on Artificial Intelligence</i> ,	953
	<i>arXiv preprint.</i>	pages 25416–25424.	954
898	Thanh Tam Nguyen, Thanh Trung Huynh, Zhao Ren,	Yue Wang, Qizhou Wang, Feng Liu, Wei Huang, Yali Du,	955
899	Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin,	Xiaojiang Du, and Bo Han. 2025b. GRU: Mitigating	956
900	and Quoc Viet Hung Nguyen. 2025. A survey of	the trade-off between unlearning and retention for	957
901	machine unlearning.	LLMs.	958
902	<i>ACM Transactions on Intelligent Systems and Technology</i> ,	<i>In Proceedings of the 42nd International</i>	959
	16(5):1–46.	<i>Conference on Machine Learning</i> , volume 267 of	960
903	Vaidehi Ramesh Patil, Peter Hase, and Mohit Bansal.	<i>Proceedings of Machine Learning Research</i> , pages	961
904	2024. Can sensitive information be deleted from	64690–64710. PMLR.	
905	LLMs? objectives for defending against extraction		
906	attacks.	Ronald L. Wasserstein and Nicole A. Lazar. 2016. The	962
907	<i>In The Twelfth International Conference on Learning Representations (ICLR 2024)</i> ,	asa statement on p-values: Context, process, and pur-	963
908	pages 45497–45514.	pose.	964
		<i>The American Statistician</i> , 70(2):129–133.	
909	Nils Reimers and Iryna Gurevych. 2019. Sentence-	James H. Wilkinson. 1963. <i>Rounding Errors in Alge-</i>	965
910	BERT: Sentence embeddings using Siamese BERT-	<i>braic Processes.</i> Prentice-Hall, Englewood Cliffs, NJ.	966
911	networks.	Reprinted by Dover, 1994.	967
912	<i>In Proceedings of the 2019 Conference on</i>		
913	<i>Empirical Methods in Natural Language Processing</i>	Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu,	968
914	<i>and the 9th International Joint Conference on Natu-</i>	Julien Demouth, and Song Han. 2023. SmoothQuant:	969
915	<i>ral Language Processing (EMNLP-IJCNLP)</i> ,	Accurate and efficient post-training quantization for	970
916	pages 3982–3992. Association for Computational Linguis-		
	tics.		

- 971 [large language models](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume
972 202 of *Proceedings of Machine Learning Research*,
973 pages 38087–38099. PMLR.
974
- 975 Zeguan Xiao, Xuanzhe Xu, Yun Chen, Yong Wang, Jian
976 Yang, Yanqing Hu, and Guanhua Chen. 2026. [Ro-](#)
977 [bust LLM unlearning against relearning attacks: The](#)
978 [minor components in representations matter](#). *arXiv*
979 *preprint*.
- 980 Jin Yao, Eli Chien, Minxin Du, Xinyao Niu, Tianhao
981 Wang, Zezhou Cheng, and Xiang Yue. 2024. [Ma-](#)
982 [chine unlearning of pre-trained large language mod-](#)
983 [els](#). In *Proceedings of the 62nd Annual Meeting of the*
984 *Association for Computational Linguistics (Volume*
985 *1: Long Papers)*, pages 8403–8419. Association for
986 Computational Linguistics.
- 987 Sangyeon Yoon, Wonje Jeung, and Albert No. 2025.
988 [R-TOFU: Unlearning in large reasoning models](#). In
989 *Proceedings of the 2025 Conference on Empirical*
990 *Methods in Natural Language Processing*, pages
991 5239–5258, Suzhou, China. Association for Com-
992 putational Linguistics.
- 993 Ashkan Yousefpour, Taeheon Kim, Ryan Sungmo Kwon,
994 Seungbeen Lee, Wonje Jeung, Seungju Han, Alvin
995 Wan, Harrison Ngan, Youngjae Yu, and Jonghyun
996 Choi. 2025. [Representation bending for large lan-](#)
997 [guage model safety](#). In *Proceedings of the 63rd*
998 *Annual Meeting of the Association for Computa-*
999 *tional Linguistics (Volume 1: Long Papers)*, pages
1000 24073–24098. Association for Computational Lin-
1001 guistics.
- 1002 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024.
1003 [Negative preference optimization: From catastrophic](#)
1004 [collapse to effective unlearning](#). In *Proceedings of*
1005 *the First Conference on Language Modeling (COLM*
1006 *2024)*.
- 1007 Zhiwei Zhang, Fali Wang, Xiaomin Li, Zongyu Wu,
1008 Xianfeng Tang, Hui Liu, Qi He, Wenpeng Yin, and
1009 Suhang Wang. 2025. [Catastrophic failure of LLM](#)
1010 [unlearning via quantization](#). In *The Thirteenth In-*
1011 *ternational Conference on Learning Representations*
1012 *(ICLR 2025)*, pages 74925–74948.

1013 Appendix

1014 This appendix provides full technical details, ex-
1015 tended analyses, and supplementary experiments.

1016 A Extended Methodology Details

1017 The core Q-ROU framework, problem formulation,
1018 and loss objectives are defined in the main text (Sec-
1019 tion 2). This section supplies supporting analytical
1020 results and the complete training algorithm.

1021 A.1 Design Rationale: Minimal Analytical 1022 Results

1023 The following statements are lightweight analytical
1024 consequences of the optimization design and evalu-
1025 ation thresholding scheme. They provide sufficient-
1026 condition diagnostics for the audited operating point
1027 and clarify which parts of the deployment story are
1028 analytical versus empirical.

1029 **Analytical scope.** Proposition 1 is a threshold-
1030 stability diagnostic under bounded log-probability
1031 perturbations, not universal robustness to all quan-
1032 tizers. Proposition 2 interprets AR as a Lagrangian
1033 relaxation of neighbor-drift control, not a global
1034 Fisher barrier. Theorem 1 explains QuantNoise as
1035 a noise-weighted Hessian-trace penalty, not a guar-
1036 antee for every production quantizer. Theorem 2
1037 bounds one matrix-subspace rotation channel un-
1038 der small updates, not semantic preservation from
1039 linear algebra alone. Theorem 4 is a Fano-style
1040 scaling statement under an explicit output-channel
1041 leakage assumption, not protection against side-
1042 channel, parameter-access, or unrestricted para-
1043 phrase attacks. Theorem 5 is a noised-model indis-
1044 tinguishability calculation under a stated retrained-
1045 distance assumption, not an empirical retraining-
1046 equivalence result.

1047 **Proposition 1 (Threshold Stability under
1048 Bounded Quantization Perturbation).** Let τ be
1049 the target-pass threshold ($\tau = 0.01$), and define
1050 the log-probability margin for a target case (x, y)
1051 as $m(x, y) = \log \tau - \log P_{\theta^*}(y|x)$. Assume quan-
1052 tization induces bounded log-probability pertur-
1053 bation $|\log P_{Q(\theta^*)}(y|x) - \log P_{\theta^*}(y|x)| \leq \epsilon_q$. If
1054 $m(x, y) > \epsilon_q$, the pass/fail decision for that case is
1055 unchanged after quantization. **Implication.** The
1056 margin term and QuantNoise jointly increase target-
1057 side separation from the threshold, reducing flip
1058 risk. *Caveat:* The bounded log-probability per-
1059 turbation assumption holds most naturally in the
1060 neighbor/general retention regime where probabil-
1061 ities are moderate. In the target regime where

$P_{\theta^*}(y|x) \ll 1$, the effective perturbation ϵ_q in log-
1062 probability space may grow—since $\log P$ becomes
1063 highly sensitive to small probability changes near
1064 zero—potentially weakening this conclusion for
1065 cases closest to threshold. 1066

**Proposition 2 (AR as Constrained Neighbor-
1067 Drift Control).** Consider the constrained problem: 1068

$$\begin{aligned} \min_{\theta} \quad & \mathcal{L}_{\text{forget}} + \lambda_r \mathcal{L}_{\text{retain}} + \mathcal{R}(\theta) \\ \text{s.t.} \quad & \mathbb{E}_{x_n} D_{\text{KL}}(P_{\theta_0}(\cdot|x_n) \| P_{\theta}(\cdot|x_n)) \leq \xi \end{aligned} \quad (5) \quad 1069$$

1070 where \mathcal{R} aggregates EWC, orthogonality, and mar-
1071 gin terms. Our penalized objective (combining
1072 Eqs. 1 and 2) is the Lagrangian relaxation of Eq. 5
1073 with multiplier λ_a .

Implication. Increasing λ_a tightens the optimiza-
1074 tion toward smaller neighbor-output drift, confirm-
1075 ing AR is a direct, tunable knob for neighbor devia-
1076 tion control. 1077

Remark (Fisher-Geometric Interpretation). 1078
1079 Under the standard small-update expansion of KL
1080 divergence around the frozen reference model, the
1081 KL constraint in Eq. 5 has a Fisher-geometric local
1082 form involving the neighbor-aggregated Fisher in-
1083 formation matrix \mathcal{I}_{NBR} : $\Delta\theta^T \mathcal{I}_{\text{NBR}} \Delta\theta$ is controlled
1084 by the active-retention multiplier and the feasible
1085 neighbor-drift budget. This encourages parame-
1086 ter updates toward the approximate null space of
1087 \mathcal{I}_{NBR} —directions that leave neighbor output dis-
1088 tributions approximately invariant. AR therefore
1089 supplies a local parameter-space interpretation of
1090 neighbor protection without requiring explicit con-
1091 cept localization (Hase et al., 2023). This Fisher
1092 null-space reading is the *primary* theoretical jus-
1093 tification for AR, but it should be read as a local
1094 approximation around the fixed reference anchor
1095 rather than as a global structural barrier. The theo-
1096 retical foundation for using Fisher information to
1097 guide selective forgetting was established by (Go-
1098 latkar et al., 2020), which formalized unlearning as
1099 achieving weight-space indistinguishability from a
1100 model retrained without the forget set. As a supple-
1101 mentary empirical tool, we additionally introduce
1102 a Wronskian-based collision-risk signal computed
1103 from short-window AR surrogate fits of hidden-
1104 state trajectories, providing lightweight online early-
1105 warning capability (detailed below). 1105

1106 A.2 Training Procedure

1107 Algorithm 1 summarizes the complete Q-ROU train-
1108 ing procedure. All trainable parameters in SLUG 1108

1109 layers are cast to FP32 for optimization stability
 1110 (even when the model backbone uses FP16), and
 1111 cast back after training.

Algorithm 1 Q-ROU Training Procedure

Require: Pre-trained model θ , forget set \mathcal{D}_f ,
 neighbor set \mathcal{D}_n , general set \mathcal{D}_g , SLUG layers
 \mathcal{S} , steps K

- 1: Compute original logits: $\{P_\theta(\cdot|x)\}_{x \in \mathcal{D}_n \cup \mathcal{D}_g}$
 - 2: Compute truth directions $\{v_l\}_{l \in \mathcal{S}}$ via PCA on TruthfulQA
 - 3: Freeze all parameters except $\{W_l\}_{l \in \mathcal{S}}$
 - 4: Cast trainable params to FP32
 - 5: **for** $k = 1, \dots, K$ **do**
 - 6: Sample batches $B_f \subset \mathcal{D}_f$, $B_n \subset \mathcal{D}_n$,
 $B_g \subset \mathcal{D}_g$
 - 7: Set $\lambda_a^{(k)} \leftarrow \lambda_a$
 - 8: **if** optional adaptive AR variant enabled
then
 - 9: Fit short-window AR responses from
 target/neighbor hidden trajectories and compute
 Wronskian risk ρ_k
 - 10: Adapt $\lambda_a^{(k)}$ based on the surrogate col-
 lision signal
 - 11: **end if**
 - 12: **if** adaptive QuantNoise variant enabled
then
 - 13: Fetch optimizer 2nd moment v_l and
 compute adaptive scale $\beta_v \propto (v_l/\bar{v})^\beta$
 - 14: Inject noise: $\tilde{W}_l \leftarrow W_l + \eta \cdot s(W_l) \cdot \beta_v$,
 $\eta \sim \mathcal{U}(-\epsilon_q, \epsilon_q)$
 - 15: **else**
 - 16: Inject standard QuantNoise into train-
 able SLUG-layer weights
 - 17: **end if**
 - 18: Compute $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{core}} + \mathcal{L}_{\text{aux}}$ using $\lambda_a^{(k)}$
 - 19: Update $\{W_l\}_{l \in \mathcal{S}}$ via Adam optimizer
 - 20: **end for**
 - 21: Cast trainable params back to original precision
 - 22: **return** modified model θ^*
-

1112 **A.3 Analytical Diagnostics for Active**
 1113 **Retention and QuantNoise**

1114 This section provides supporting analytical diagnos-
 1115 tics for the two core optimization mechanisms in
 1116 Q-ROU: QuantNoise training and Active Retention.
 1117 The goal is to clarify why these mechanisms are
 1118 plausible under explicit assumptions, not to present
 1119 a stand-alone theory of irreversible unlearning.

A.3.1 QuantNoise and its Adaptive Variant 1120
Induce Flat Maxima 1121

Let $w \in \mathbb{R}^n$ be the parameters of the SLUG lay- 1122
 ers. The standard QuantNoise mechanism injects 1123
 uniform noise $\eta_i \sim \mathcal{U}(-\epsilon_0, \epsilon_0)$ into the weights 1124
 during each forward pass. The Adaptive Quant- 1125
 Noise extension directionally scales the perturba- 1126
 tion bounds via $\epsilon_{q,i} = \epsilon_0 \beta_{v,i}$, where $\beta_{v,i} \propto (v_i/\bar{v})^\gamma$ 1127
 is derived from the Adam optimizer’s second mo- 1128
 ment v and $\gamma > 0$ is a tunable exponent. Con- 1129
 ceptually, this does not freeze the model; rather, it 1130
 broadens the local optimum so that quantization per- 1131
 turbations are absorbed within a flatter basin. This 1132
 design philosophy—regularizing during training to 1133
 facilitate post-hoc modifications—shares concep- 1134
 tual foundations with (Thudi et al., 2022), which 1135
 proposed standard deviation (SD) loss to constrain 1136
 weight trajectories and reduce verification error 1137
 during unlearning. While SD loss targets weight- 1138
 space proximity to retrained models, QuantNoise 1139
 targets loss-landscape flatness for quantization ro- 1140
 bustness. The importance of activation-aware quan- 1141
 tization strategies is further emphasized by (Xiao 1142
 et al., 2023), which demonstrated that activation 1143
 outliers are the primary bottleneck for low-bit quan- 1144
 tization in LLMs, motivating our focus on loss- 1145
 landscape geometry rather than weight magnitude 1146
 alone. Comprehensive surveys of quantization tech- 1147
 niques (Nagel et al., 2021) establish the theoretical 1148
 foundations for post-training quantization (PTQ) 1149
 and quantization-aware training (QAT), providing 1150
 the broader context within which QuantNoise op- 1151
 erates as a specialized regularization strategy for 1152
 unlearning robustness. 1153

Theorem 1 (QuantNoise Regularizes a 1154
Weighted Hessian Trace). Let $\mathcal{L}(w)$ be the 1155
 loss function, assumed four times continuously 1156
 differentiable. Let $\eta \in \mathbb{R}^n$ with independent com- 1157
 ponents $\eta_i \sim \mathcal{U}(-\epsilon_{q,i}, \epsilon_{q,i})$, where $\epsilon_{q,i} = \epsilon_0 \beta_{v,i}$ 1158
 (with $\beta_{v,i} = 1$ for standard QuantNoise). Then the 1159
 expected perturbed loss satisfies: 1160

$$\mathbb{E}_\eta[\mathcal{L}(w + \eta)] = \mathcal{L}(w) + \frac{\epsilon_0^2}{6} \sum_{i=1}^n \beta_{v,i}^2 \nabla_{ii}^2 \mathcal{L}(w) + \mathcal{O}(\epsilon_0^4) \quad (6) \quad 1161$$

Proof. We perform a second-order Taylor expan- 1162
 sion of $\mathcal{L}(w + \eta)$ around w : 1163

$$\mathcal{L}(w + \eta) = \mathcal{L}(w) + \nabla \mathcal{L}(w)^\top \eta + \frac{1}{2} \eta^\top H \eta + \mathcal{O}(\|\eta\|^3) \quad (7) \quad 1164$$

1165 where $H = \nabla^2 \mathcal{L}(w)$ is the Hessian matrix. Taking
 1166 the expectation over η :

$$\begin{aligned} \mathbb{E}_\eta[\mathcal{L}(w + \eta)] &= \mathcal{L}(w) + \nabla \mathcal{L}(w)^\top \mathbb{E}[\eta] \\ &\quad + \frac{1}{2} \sum_{i,j} H_{ij} \mathbb{E}[\eta_i \eta_j] + \mathcal{O}(\epsilon_0^4) \end{aligned} \quad (8)$$

1167 Since η_i are independent uniform variables with
 1168 mean 0, we have $\mathbb{E}[\eta_i] = 0$ and $\mathbb{E}[\eta_i \eta_j] = 0$ for
 1169 $i \neq j$. The second moment of $\mathcal{U}(-\epsilon_{q,i}, \epsilon_{q,i})$ is
 1170 $\text{Var}(\eta_i) = (2\epsilon_{q,i})^2/12 = \epsilon_{q,i}^2/3 = \epsilon_0^2 \beta_{v,i}^2/3$. Thus:

$$\begin{aligned} \frac{1}{2} \sum_{i,j} H_{ij} \mathbb{E}[\eta_i \eta_j] &= \frac{1}{2} \sum_{i=1}^n H_{ii} \cdot \frac{\epsilon_0^2 \beta_{v,i}^2}{3} \\ &= \frac{\epsilon_0^2}{6} \sum_{i=1}^n \beta_{v,i}^2 H_{ii} \end{aligned} \quad (9)$$

1173 Substituting back, and noting that the third moment
 1174 $\mathbb{E}[\eta_i^3]$ vanishes by the symmetry of the uniform
 1175 distribution (leaving the fourth-order term as the
 1176 leading nonzero remainder), yields the theorem. ■

1177 Equivalently, if $\Sigma_\eta = \text{diag}(\epsilon_{q,i}^2/3)$, the second-
 1178 order penalty is $\frac{1}{2} \text{tr}(H \Sigma_\eta)$. Thus minimizing the
 1179 QuantNoise-perturbed loss jointly minimizes the
 1180 original loss and a *noise-weighted* Hessian trace. In
 1181 the standard isotropic formulation ($\beta_{v,i} = 1$), this
 1182 reduces to an unweighted trace penalty up to the
 1183 constant $\epsilon_0^2/6$. In the adaptive variant, because $\beta_{v,i}^2$
 1184 is proportional to the optimizer’s second moment
 1185 (which estimates the uncentered gradient variance
 1186 and is often correlated with high curvature), the
 1187 penalty is concentrated on high-variance directions.
 1188 This is a design rationale for low-bit stability, not
 1189 a standalone guarantee that every production quan-
 1190 tizer will preserve the same operating point.

1191 A.3.2 Active Retention: Surrogate-Based 1192 Collision-Risk Monitoring

1193 The theoretical treatment of Active Retention pro-
 1194 ceeds in two stages. First, we use the orthogonal-
 1195 ity constraint to bound selected matrix-subspace
 1196 rotation via Davis-Kahan (Theorem 2), providing
 1197 one structural diagnostic for why localized edits
 1198 reduce collateral drift. Second, we show that the
 1199 KL penalty confines parameter updates to the ap-
 1200 proximate null space of the neighbor-aggregated
 1201 Fisher information \mathcal{I}_{NBR} (Proposition 2), provid-
 1202 ing an assumption-light design consequence for
 1203 neighbor preservation. As a complementary em-
 1204 pirical tool, the Wronskian risk signal—computed
 1205 from short-window AR surrogate fits of hidden-
 1206 state trajectories—serves as a lightweight online

early-warning indicator for concept collision in non- 1207
 catastrophic update regimes where the AR surro- 1208
 gate fit remains valid, motivating an optional adap- 1209
 tive variant (Observation 1, Motivation 1). 1210

Theoretical hierarchy and scope. The primary 1211
 justification for Active Retention is the assumption- 1212
 light local result in Proposition 2 and the Fisher- 1213
 geometric remark (Section A): under a fixed refer- 1214
 ence anchor, small updates, and the quadratic KL 1215
 approximation, the KL penalty biases updates to- 1216
 ward the approximate null space of \mathcal{I}_{NBR} . As an 1217
 optional supplement, when hidden-state trajectories 1218
 are locally approximated by short-window AR(p) 1219
 surrogate models (validity conditions stated below), 1220
 the Wronskian provides a frequency-domain early- 1221
 warning signal for concept collision. The theo- 1222
 retical hierarchy thus consists of two levels: (i) 1223
 the primary **Davis-Kahan + Fisher null-space** 1224
 results (Theorem 2, Proposition 2), which bound 1225
 knowledge-subspace rotation and neighbor drift 1226
 under local approximation assumptions; and (ii) 1227
 the supplementary **Wronskian** ρ_k (Observation 1), 1228
 which provides a practical monitoring tool for im- 1229
 minent root collision under the AR surrogate as- 1230
 sumption. The Wronskian heuristic therefore com- 1231
 plements these analytical results as a practical mon- 1232
 itoring tool rather than replacing them. 1233

1234 A.4 Protection of Knowledge Subspaces via 1235 Orthogonality Constraint (Application of 1236 Davis-Kahan Theorem)

Q-ROU modifies the weight matrix W of SLUG 1237
 layers into $W + \Delta W$ to suppress target responses. 1238
 However, these weight matrices also encode gener- 1239
 al knowledge (e.g., scientific, geographical, and 1240
 mathematical facts) completely unrelated to the tar- 1241
 get domains. Here, we use the Davis-Kahan sin Θ 1242
 theorem (Davis and Kahan, 1970) as a local per- 1243
 turbation diagnostic: when the SLUG-layer update 1244
 is small relative to the relevant spectral gap, the 1245
 dominant right-singular subspace cannot rotate ar- 1246
 bitrarily. This is not a proof of semantic deletion or 1247
 permanent preservation; it is a sufficient-condition 1248
 bound on one matrix-level source of collateral drift. 1249

**Theorem 2 (Protection of Knowledge Sub- 1250
 space via Orthogonality Constraint).** Let $\sigma_1 \geq$ 1251
 $\sigma_2 \geq \dots \geq 0$ be the singular values of the weight 1252
 matrix $W \in \mathbb{R}^{m \times n}$ of a SLUG layer. Let \mathcal{V} be the 1253
 subspace spanned by the top- k right singular vec- 1254
 tors of W , and let $\tilde{\mathcal{V}}$ be the corresponding subspace 1255
 of the post-unlearning matrix $\tilde{W} = W + \Delta W$. As- 1256

1257 suming $\delta = \sigma_k^2 - \sigma_{k+1}^2 > 0$, the following bound
 1258 holds:

$$1259 \quad \|\sin \Theta\| \leq \frac{d(2a + d)}{\delta} \quad (10)$$

1260 where $a := \|W\|_{\text{op}}$, $d := \|\Delta W\|_{\text{op}}$, and $\delta := \sigma_k^2 -$
 1261 σ_{k+1}^2 , with $\Theta = \Theta(\mathcal{V}, \tilde{\mathcal{V}})$.

1262 *Proof.* We perform Singular Value Decom-
 1263 position (SVD) on W and denote the subspace
 1264 spanned by the top- k right singular vectors as
 1265 $\mathcal{V} = \text{span}\{v_1, \dots, v_k\}$. To apply the Davis-Kahan
 1266 theorem, we construct the symmetric matrix $M =$
 1267 $W^T W$. The eigenvectors of M coincide with the
 1268 right singular vectors of W , associating with eigen-
 1269 values σ_i^2 . The post-unlearning matrix takes the
 1270 form $\tilde{M} = \tilde{W}^T \tilde{W} = W^T W + E$, where the perturba-
 1271 tion matrix E is given by $E = W^T \Delta W + \Delta W^T W +$
 1272 $\Delta W^T \Delta W$. By applying the triangle inequality, the
 1273 operator norm of E is bounded from above:

$$1274 \quad \|E\|_{\text{op}} \leq 2\|W\|_{\text{op}}\|\Delta W\|_{\text{op}} + \|\Delta W\|_{\text{op}}^2 \quad (11)$$

1275 Given the spectral gap of M is $\delta = \sigma_k^2 - \sigma_{k+1}^2$,
 1276 applying the Davis-Kahan $\sin \Theta$ theorem (Davis
 1277 and Kahan, 1970) immediately yields the bound
 1278 stated in the theorem. (Note: the bound is non-
 1279 vacuous only when $\|E\|_{\text{op}} < \delta$; under Q-ROU’s
 1280 short-step regime with $\|\Delta W\|_{\text{op}}/\|W\|_{\text{op}} < 1\%$, this
 1281 condition is comfortably satisfied in practice.) ■

1282 **Implications and Numerical Validation.** This
 1283 bound shows that the rotation of the selected matrix
 1284 subspace is restricted by the magnitude of pertur-
 1285 bation $\|\Delta W\|_{\text{op}}$ and inversely proportional to the
 1286 spectral gap δ . We use it as one diagnostic explain-
 1287 ing why localized updates can be less damaging
 1288 than unconstrained full-layer edits. First, the use of
 1289 **localized updates** via SLUG over a minimal num-
 1290 ber of steps ensures that the relative perturbation
 1291 ratio $\|\Delta W\|_{\text{op}}/\|W\|_{\text{op}}$ remains extremely small (typ-
 1292 ically $< 1\%$). Second, the activation-based SLUG
 1293 profiling inherently selects layers with **intrinsic**
 1294 **gap preservation**, tending toward those with mod-
 1295 erately large spectral gaps which naturally fortifies
 1296 the stability of the bound. Third, the **orthogonal**
 1297 **projection** through the $\mathcal{L}_{\text{ortho}}$ constraint and PGD
 1298 projection reduces update components aligned with
 1299 the truth direction v_l . Since truth directions are
 1300 correlated with broad factual directions, this helps
 1301 limit one class of collateral update directions while
 1302 restricting the principal rotation $\|\sin \Theta\|$.

1303 Numerical simulations across varying matrix
 1304 scales ($m = 32 \sim 3072$) confirm that the algebraic
 1305 bound is satisfied in the tested perturbation regime

(typical tightness ratio ~ 0.13). Under the prevail- 1306
 ing operating regime of Q-ROU, the corresponding 1307
 matrix-subspace rotation bound is small (approximate 1308
 bounding angle $\leq 1.5^\circ \sim 8^\circ$). We do not equate 1309
 this matrix bound with a semantic guarantee; the 1310
 semantic claim is evaluated empirically through 1311
 neighbor, generation, and extraction audits. 1312

Remark 2 (Spectral Gaps and Representa- 1313
tional Geometry). The central role of the spec- 1314
 tral gap δ in our bound has a natural parallel in 1315
 the representation geometry literature. Karkada 1316
et al. (Karkada et al., 2026) show that when 1317
 co-occurrence statistics between semantically rel- 1318
 ated words exhibit translation symmetry, learned 1319
 embeddings organize into Fourier modes whose 1320
 eigenvalues—and thus spectral gaps—grow with 1321
 the number of vocabulary items sharing the same 1322
 latent variable (a “collective effect”). Applying 1323
 the Davis-Kahan theorem to this setting, they prove 1324
 that such representations are robust to localized 1325
 perturbation of co-occurrence entries, precisely be- 1326
 cause the spectral gaps overwhelm the perturbation 1327
 magnitude. This parallel suggests that the spectral 1328
 gaps protecting knowledge subspaces in our The- 1329
 orem 2 may not be incidental structural features 1330
 of pre-trained weight matrices, but rather reflect 1331
 the systematic organization of knowledge repre- 1332
 sentations around shared latent variables. In ei- 1333
 ther case—whether spectral gaps arise from collec- 1334
 tive co-occurrence structure or from the intrinsic 1335
 spectrum of pre-trained weights—the Davis-Kahan 1336
 mechanism provides a unified explanation for why 1337
 localized perturbations (be they statistical noise or 1338
 unlearning updates) preserve the dominant repre- 1339
 sentational structure. 1340

The theoretical foundation for using Fisher infor- 1341
 mation to guide selective forgetting was established 1342
 by (Golatkar et al., 2020), which formalized un- 1343
 learning as achieving weight-space indistinguishability 1344
 from a model retrained without the forget set. 1345
 Our Fisher null-space argument (Proposition 2) ex- 1346
 tends this framework to the neighbor-preservation 1347
 setting, where the goal is to confine updates to di- 1348
 rections that leave neighbor output distributions 1349
 approximately invariant. ◊ 1350

Concept entanglement—where suppressing a tar- 1351
 get concept also damages semantically adjacent 1352
 neighbor concepts—is a recurring failure mode of 1353
 standard gradient unlearning (Wang et al., 2026). 1354
 Our main analytical protection is the Fisher null- 1355
 space interpretation of AR (Proposition 2) together 1356
 with the orthogonality regularizer (Theorem 2). 1357

1358 As a supplementary practical tool, we introduce
 1359 a Wronskian-based warning signal that monitors
 1360 collision risk at the surrogate-model level. This
 1361 signal is computed by fitting short-window $\text{AR}(p)$
 1362 surrogate models to hidden-state trajectories across
 1363 the token dimension and tracking whether their char-
 1364 acteristic polynomial roots converge.

1365 **Modeling assumption (locally quasi-**
 1366 **stationary AR).** We do not claim that hidden state
 1367 trajectories are generated by a stationary $\text{AR}(p)$
 1368 process in the strict statistical sense. Rather, we
 1369 model the token-dimension evolution of a concept’s
 1370 hidden states as a *locally quasi-stationary* $\text{AR}(p)$
 1371 process—analogue to linear predictive coding
 1372 (LPC) in speech processing, where a highly
 1373 non-linear, time-variant signal is approximated
 1374 by short-window stationary AR segments. The
 1375 knowledge structure of concept A is encoded in its
 1376 characteristic polynomial $\phi_A(z) = \sum_{k=0}^p \alpha_k z^{p-k}$,
 1377 where z is the complex shift operator.

1378 This modeling framework is meaningful when
 1379 the following empirically verifiable conditions
 1380 hold:

- 1381 1. **AR fit quality:** the short-window $\text{AR}(p)$ fit
 1382 to hidden-state trajectories achieves sufficient
 1383 goodness-of-fit ($R^2 \gtrsim 0.8$), indicating that the
 1384 quasi-stationary approximation captures the
 1385 dominant dynamics;
- 1386 2. **Sufficient token length:** concept-specific
 1387 texts are long enough ($T \gg p$, typically
 1388 $T > 50$ tokens) to yield stable AR coefficient
 1389 estimates;
- 1390 3. **Initial root separation:** the baseline (pre-
 1391 unlearning) characteristic polynomials ϕ_A, ϕ_B
 1392 have well-separated roots ($\delta_0 > 0$), so that
 1393 the surrogate models start in a non-degenerate
 1394 regime.

1395 When these conditions are reasonably satisfied,
 1396 the Wronskian risk signal provides useful early-
 1397 warning information about potential concept colli-
 1398 sion. When they are violated (e.g., very short texts
 1399 or poor AR fit), the signal should be treated with
 1400 caution. In all cases, the primary analytical results
 1401 (Proposition 2, Theorem 2) remain the main design
 1402 reference regardless of AR surrogate fit quality.

1403 Why Wronskian as the early-warning signal?

1404 When hidden-state trajectories of two concepts A
 1405 and B are locally approximated by $\text{AR}(p)$ surro-
 1406 gate models, each concept is characterized by its

surrogate characteristic polynomial $\phi_A(z), \phi_B(z)$. 1407
 Root sharing ($\phi_A(\zeta^*) = \phi_B(\zeta^*) = 0$) indicates 1408
 that both surrogate models resonate at the same 1409
 frequency—a condition correlated with concept en- 1410
 tanglement in practice. However, repeatedly esti- 1411
 mating full polynomial resultants online is numer- 1412
 ically brittle under noisy short-window AR esti- 1413
 mates. The Wronskian provides a computationally 1414
 tractable, gradient-sensitive proxy for detecting im- 1415
 minent root collision, without requiring full polyno- 1416
 mial resultant computation. It is (i) computable on- 1417
 line during training, (ii) sensitive to *imminent* root 1418
 collision rather than only post-hoc drift, and (iii) 1419
 more informative than generic embedding-space 1420
 distances (cosine/CKA/mean-feature gaps), which 1421
 cannot test whether two dynamics are approaching 1422
 root-sharing. 1423

Lemma 1 (Local root-gap sensitivity of Wron-
skian evaluated on the unit circle). Let $\phi_A(z)$ and 1424
 $\phi_B(z)$ be monic polynomials of degree p represent- 1425
 ing the characteristic responses of concepts A and 1426
 B . To analyze these responses continuously in the 1427
 frequency domain, we evaluate the polynomials on 1428
 the unit circle $z = e^{i\omega}$ for $\omega \in [0, 2\pi)$. The con- 1429
 tinuous Wronskian of these frequency responses is 1430
 defined as: 1431

$$1432 \quad W(\phi_A, \phi_B)(\omega) = \phi_A(e^{i\omega}) \frac{d}{d\omega} \phi_B(e^{i\omega}) \quad (12) \quad 1433$$

$$1434 \quad - \phi_B(e^{i\omega}) \frac{d}{d\omega} \phi_A(e^{i\omega}) \quad 1435$$

When evaluating at a shared simple root $\zeta^* =$ 1434
 $e^{i\omega_0}$, $W(\phi_A, \phi_B)(\omega_0) = 0$. Furthermore, for close 1435
 roots ζ_A, ζ_B near ω_0 , the magnitude of the Wron- 1436
 skian serves as a first-order proxy for the root gap: 1437
 $|W(\omega_0)| = \Theta(|\zeta_A - \zeta_B|)$ and $W \rightarrow 0$ as roots col- 1438
 lide. 1439

Proof sketch. Let ω_A, ω_B be the frequencies cor- 1440
 responding to nearby roots $\zeta_A = e^{i\omega_A}, \zeta_B = e^{i\omega_B}$. 1441
 Near $\omega_0 \approx \omega_A \approx \omega_B$, expand $\phi_A(e^{i\omega_0}) = c_A(\omega_0 -$ 1442
 $\omega_A) + \mathcal{O}((\omega_0 - \omega_A)^2)$ and $\phi_B(e^{i\omega_0}) = c_B(\omega_0 -$ 1443
 $\omega_B) + \mathcal{O}((\omega_0 - \omega_B)^2)$, where $c_A = \phi'_A(e^{i\omega_A}) \cdot i e^{i\omega_A}$ 1444
 and similarly for c_B . The derivatives $\frac{d}{d\omega} \phi_A(e^{i\omega_0})$ 1445
 evaluate to $c_A + \mathcal{O}(\omega_0 - \omega_A)$. Substituting into 1446
 the Wronskian determinant and retaining leading- 1447
 order terms, the dominant contribution is propor- 1448
 tional to $c_A c_B (\omega_A - \omega_B) + \mathcal{O}((\omega_A - \omega_B)^2)$. Since 1449
 $|\omega_A - \omega_B| = \Theta(|\zeta_A - \zeta_B|)$ for roots on the unit 1450
 circle, the claim follows. \square 1451

Implication. Monitoring $\rho_k = \min_{\omega \in \Omega} |W_k(\omega)|$ 1452
 provides a low-overhead collision-risk signal that is 1453
 more operationally useful than generic embedding 1454

1455 distances for detecting imminent root collision in
1456 the surrogate models.

1457 **Surrogate Collision Detection.** When unlearn-
1458 ing a target concept A (causing its surrogate charac-
1459 teristic roots to shift), a neighbor concept B faces
1460 collision risk if the parameter updates cause the
1461 surrogate polynomials $\phi_A(z)$ and $\phi_B(z)$ to drift
1462 toward sharing a root. To detect this impending
1463 collision, Q-ROU’s optional adaptive variant mon-
1464 itors the surrogate Wronskian risk score $\rho_k =$
1465 $\min_{\omega \in \Omega} |W_k(\omega)|$.

1466 This continuous monitoring is meaningful be-
1467 cause, under the AR surrogate model, roots vary
1468 continuously with the fitted coefficients away from
1469 degenerate cases. The argument relies on three lo-
1470 cal conditions: (a) polynomial roots are continuous
1471 functions of their coefficients (Ostrowski’s classical
1472 theorem); (b) the AR(p) coefficients are continuous
1473 functions of the model parameters θ via the Yule-
1474 Walker equations $a = \Gamma_p^{-1} \gamma$ as long as Γ_p remains
1475 non-singular; and (c) gradient-based training moves
1476 θ in small steps. Together, (a)–(c) make abrupt un-
1477 noticed root collisions unlikely within the surrogate
1478 model, motivating ρ_k as an early-warning signal. \diamond

1479 *Intuition.* The Wronskian is an online collision
1480 alarm at the surrogate level: when two concepts’
1481 AR surrogate models start sharing roots, the signal
1482 drops toward zero before downstream retention met-
1483 rics collapse. **Observation 1 (Wronskian Drop**
1484 **Signals Surrogate Root Collision).** Under the AR
1485 surrogate model, if the roots of ϕ_A and ϕ_B converge
1486 (i.e., $\Delta\zeta \rightarrow 0$), the Wronskian evaluated on the unit
1487 circle satisfies $W(\phi_A, \phi_B)(\omega) \rightarrow 0$.

1488 *Proof.* By Lemma 1, for nearby simple roots
1489 $\zeta_A = e^{i\omega_A}$, $\zeta_B = e^{i\omega_B}$ near ω_0 , $|W(\omega_0)| =$
1490 $\Theta(|\zeta_A - \zeta_B|)$. Therefore $\Delta\zeta \rightarrow 0$ implies $W(\omega_0) \rightarrow$
1491 0 and hence $\rho_k = \min_{\omega} |W_k(\omega)| \rightarrow 0$. At exact
1492 root sharing ($\phi_A(\zeta^*) = \phi_B(\zeta^*) = 0$), the solu-
1493 tions become linearly dependent and $W = 0$ exactly.
1494 Thus, a decreasing Wronskian acts as an early warn-
1495 ing for root collision in the surrogate models. ■

1496 *Scope.* This observation holds at the surrogate-
1497 model level: it characterizes the behavior of the
1498 AR(p) approximation to hidden-state trajectories.
1499 Whether surrogate root collision reliably predicts
1500 actual concept entanglement in the full non-linear
1501 neural model depends on the fidelity of the AR
1502 surrogate fit. We therefore treat the Wronskian
1503 as an optional monitoring heuristic rather than a
1504 required component of Q-ROU’s core analytical
1505 results (Proposition 2, Theorem 2).

In implementation this is realized as a monotone
1506 barrier controller: 1507

$$\lambda_a^{(k+1)} = \begin{cases} \min(\lambda_a^{(k)} + \Delta\lambda, \lambda_{\max}), & \rho_k < \tau_W, \\ \lambda_a^{(k)}, & \text{otherwise,} \end{cases} \quad (13) \quad 1508$$

1509 where $\rho_k = \min_{\omega \in \Omega} |W_k(\omega)|$. This directly maps
1510 low Wronskian events to tighter neighbor-drift con-
1511 straints.

1512 **Motivation 1 (Empirical Motivation: Adap-**
1513 **tive AR Variant Resists Concept Collision).** 1514
1515 Proposition 2 establishes that a sufficiently large
1516 fixed λ_a confines parameter updates to the approxi-
1517 mate null space of \mathcal{I}_{NBR} , protecting neighbor knowl-
1518 edge. The adaptive variant goes further: by monitor-
1519 ing the surrogate Wronskian risk signal ρ_k (Observ-
1520 ation 1), it detects *when* collision risk is elevated in
1521 the surrogate model and dynamically increases $\lambda_a^{(k)}$
1522 to restore the safety margin. When the adaptive vari-
1523 ant of Q-ROU detects $W_k < \tau_W$ and boosts the Ac-
1524 tive Retention penalty $\lambda_a^{(k)}$, the parameter step size
1525 is throttled such that $\|\Delta\theta_B\|_F \leq \mathcal{O}(1/\lambda_a^{(k)})$. Under
1526 the local Fisher approximation (Proposition 2), stan-
1527 dard polynomial root perturbation theory (Wilkin-
1528 son, 1963) then bounds root displacement in the
1529 surrogate model, providing an empirical motivation
1530 for the ability of this adaptive variant to prevent con-
1531 cept entanglement.

1532 *Algebraic Argument.* Considering the local
1533 stationary conditions of optimization, $\nabla_{\theta_B} \mathcal{L}_f +$
1534 $\lambda_a \mathcal{I}(\theta_B) \Delta\theta_B \approx 0$, the update step relates to the
1535 gradient as $\Delta\theta_B \approx -\frac{1}{\lambda_a} \mathcal{I}(\theta_B)^{-1} \nabla_{\theta_B} \mathcal{L}_f$. Under
1536 bounded local curvature, this gives the scaling
1537 $\|\Delta\theta_B\|_F = \mathcal{O}(1/\lambda_a)$. This limits the change in
1538 the surrogate polynomial coefficients of B , which
1539 by standard perturbation theory for polynomial
1540 roots (Wilkinson, 1963) suggests a drift scale
1541 $|\Delta\zeta_i^{(B)}| \leq K(\lambda_a)^{-1}$, where K depends on the poly-
1542 nomial condition number. If the maximum root
1543 displacement is below the minimal separating dis-
1544 tance $\delta_0/2$ to the approaching roots of ϕ_A , no root
1545 crossings occur in the surrogate model. Therefore,
1546 dynamically increasing $\lambda_a^{(k)}$ when the Wronskian
1547 drops is a plausible barrier strategy for maintaining
1548 distinct surrogate roots.

1549 This indicates that the Wronskian-Adaptive vari-
1550 ant provides a practical safeguard when surrogate
1551 root collision is imminent. We acknowledge cer-
1552 tain limitations: identifying exact Lipschitz con-
1553 stants for neural dynamics is practically difficult,
1554 the Fisher Information Matrix provides a 2nd-order
1555 local approximation that accumulates error over

multiple training steps, and the AR surrogate fit is itself an approximation. Thus, this derivation serves as a motivation for using the Wronskian signal as optional online instrumentation rather than as a globally strict formal claim. In the standard Q-ROU formulation, a sufficiently large fixed λ_a provides empirical safe separation globally, albeit with less dynamic responsiveness.

A.4.1 Quantization Robustness: From Constrained Mobility to Flat Basins

Remark 1 (Asymmetric Quantization Sensitivity under AR Constraint). Let θ^* be the post-unlearning parameters obtained under Active Retention with coefficient λ_a . Let $Q(\theta^*)$ denote the INT4-quantized parameters, and define the per-case *quantization shift* $\delta(x, y) = P_{Q(\theta^*)}(y|x) - P_{\theta^*}(y|x)$. AR constrains neighbor probabilities to the high-confidence regime $P_{\theta^*}(y|x_n) \gg \tau$ where the softmax Jacobian is near-singular, while target probabilities are pushed to the low-probability regime $P_{\theta^*}(y|x_t) \approx \tau$. In this configuration, the *sensitivity* of output probabilities to weight perturbation is strongly asymmetric:

$$\left| \frac{\partial P}{\partial w} \right|_{x_n} \ll \left| \frac{\partial P}{\partial w} \right|_{x_t} \quad (14)$$

Analysis. The softmax function $\sigma(z)_k = e^{z_k} / \sum_j e^{z_j}$ has sensitivity $\partial \sigma_k / \partial z_i$ that is maximal at intermediate probabilities ($P \approx 0.5$) and minimal at extremes ($P \approx 0$ or $P \approx 1$). AR pushes neighbor probabilities into the saturated high-confidence regime, making them insensitive to small weight perturbations from quantization. Target probabilities, however, are *highly sensitive* to perturbations due to the steep gradient in the low-probability regime.

A key observation here is that this sensitivity is *not* directionally biased at first order: symmetric weight perturbations from INT4 rounding produce zero-mean logit changes, and hence zero-mean probability changes to leading order. At second order, Jensen’s inequality might suggest that the convexity of e^{z_k} in the numerator would increase the target probability; however, in the full softmax ratio $P_k = e^{z_k} / \sum_j e^{z_j}$, the same convexity effect applies to the denominator, and these contributions approximately cancel for small perturbations. The net directional bias of quantization perturbation on the softmax output is therefore *not* determined by simple single-variable convexity arguments, but depends on the global structure of the loss landscape

around the optimized point. In AR-equipped models, we empirically observe $\Delta_{\text{zombie}} < 0$ (Table 34), consistent with the composite KL-to-uniform objective reshaping the local curvature asymmetrically such that INT4 rounding tends to nudge borderline targets further below threshold. The full nonlinear dynamics of multi-layer transformer quantization involve higher-order interactions between layers that are not captured by single-layer softmax analysis.

Practical consequence. The key empirical finding remains robust: neighbor probabilities are effectively anchored by AR against quantization perturbation, while target probabilities—being in a highly sensitive regime—experience measurable shifts. The observed direction of these shifts (synergistic rather than antagonistic) is a favorable property of the Q-ROU loss landscape geometry, confirmed deterministically across three random seeds.

◇ *Scope caveat.* This analysis explains the *asymmetric sensitivity* (why targets are affected more than neighbors under quantization) rigorously, while the *directional* component (why shifts tend to be synergistic) is an empirical regularity supported by the experimental evidence (+0 to +1 target cases under INT4 in AR-equipped baselines at 40 steps) rather than a first-principles prediction.

Constrained-mobility regime and sign persistence. A complementary perspective on why short-step unlearning edits are vulnerable to quantization comes from the *constrained-mobility* of weight updates. In the sub-bit model compression literature, Sakai and Ichikawa (Sakai and Ichikawa, 2026) formalized a *sign lock-in* phenomenon: under standard SGD training, weight signs (\pm) initialized randomly are rarely flipped, because sign changes require the trajectory to cross through a narrow boundary neighborhood near zero—an exponentially unlikely event under bounded updates. Their stopping-time analysis shows that the number of effective sign flips follows a geometric tail, and this persistence strengthens with model scale.

We invoke sign lock-in as a *formally motivated geometric prior*. Although Sakai and Ichikawa’s (Sakai and Ichikawa, 2026) proof is stated for scheduled SGD, Remark D.11 of their work establishes that the geometric-tail conclusion transfers to any optimizer satisfying a bounded-update condition, including AdamW, provided that the effective per-step increment $\|p_{t+1}\|_\infty$ remains bounded on a high-probability event. In our short-

1655 step (≤ 40 steps) unlearning regime with gradi- 1704
1656 ent clipping, this condition is empirically satisfac- 1705
1657 ed. This provides intuition for why, in our un- 1706
1658 learning regime, many edited weights do not move 1707
1659 far enough to stably cross INT4 quantization-cell 1708
1660 boundaries. When post-training rounding maps 1709
1661 these borderline edits back toward their pre-edit 1710
1662 quantization cells, the Zombie Delta phenomenon 1711
1663 arises—targeted knowledge partially resurfaces un- 1712
1664 der quantization. 1713

1665 **Flat basin strategy.** Forcing larger updates to 1714
1666 cross more quantization boundaries could reduce 1715
1667 this failure mode, but at the cost of increased risk 1716
1668 of collateral neighbor damage (Neighbor Collapse). 1717
1669 Our strategy instead achieves robustness *without* re- 1718
1670 lying on aggressive displacement: Adaptive Quant- 1719
1671 Noise regularizes the Hessian curvature (Theorem 1720
1672 1), guiding optimization toward flat basins whose 1721
1673 plateaus are wider than the expected INT4 round- 1722
1674 ing perturbation radius. This approach is partic- 1723
1675 ularly critical given the emergence of outlier fea- 1724
1676 tures in large-scale transformers (Dettmers et al., 1725
1677 2022), where specific hidden dimensions exhibit ex- 1726
1678 treme magnitudes that dominate quantization error. 1727
1679 While LLM.int8() addresses outliers through mixed- 1728
1680 precision decomposition during inference, Quant- 1729
1681 Noise preemptively flattens the loss landscape dur- 1730
1682 ing unlearning, ensuring that even outlier-sensitive 1731
1683 dimensions remain robust to INT4 rounding. 1732

1684 Consequently, even when individual weight up- 1733
1685 dates remain modest, the model can maintain res- 1734
1686 ilience after quantization: the flat basin absorbs 1735
1687 rounding errors rather than allowing them to undo 1736
1688 the forget edit. This provides an empirical recon- 1737
1689 ciliation between the constrained-mobility regime 1738
1690 and quantization robustness, confirmed determinis- 1739
1691 tically across three random seeds with $\Delta_{\text{zombie}} \approx 0$ 1740
1692 for the full Q-ROU framework. 1741

1693 B Extended Experimental Setup 1742

1694 The core experimental settings are described in Sec- 1743
1695 tion 3 of the main text. This section provides ex- 1744
1696 tended design notes. 1745

1697 **Design Note on Forget Objectives.** Q-ROU uses 1746
1698 a bounded KL-to-uniform forget objective (defined 1747
1699 in Section 2), whereas the primary baselines use 1748
1700 unbounded gradient ascent on the forget set. This 1749
1701 design choice is deliberate: unbounded loss max- 1750
1702 imization is known to cause numerical instability 1751
1703 and model degradation (Jang et al., 2023), as evi-

1704 denced by the degenerate generation observed for 1705
1706 GA variants under extended training (Section D.4). 1707
1708 To control for this confound, our AR transplant ex- 1709
1710 periments (Section C.5) add AR *to the baselines'* 1711
1712 *own objectives*—that is, GA+AR retains gradient 1713
1714 ascent as its forget mechanism—thereby isolating 1714
1715 the contribution of Active Retention from the choice 1716
1716 of forget-loss formulation. The fact that GA+AR 1717
1717 recovers 8/8 neighbor retention (vs. 0/8 without 1718
1718 AR) while still using unbounded gradient ascent 1719
1719 demonstrates that AR’s benefit is independent of 1720
1720 the forget objective. We further address this distinc- 1721
1721 tion in Section D.4, where a 3×3 hyperparameter 1722
1722 sweep over learning rate and AR weight confirms 1723
1723 that the degeneration–forgetting tradeoff is intrinsic 1724
1724 to gradient ascent rather than an artifact of specific 1725
1725 hyperparameter choices. 1726

1727 **Hyperparameter Fairness in Multi-Entity Set-** 1728
1728 **tings.** Baseline learning rates (3×10^{-4} for 0.5B, 1729
1729 1×10^{-4} for 3B) were selected in the single- 1730
1730 entity setting and applied without modification to 1731
1731 multi-entity experiments to ensure fair comparison. 1732
1732 While the multi-entity regime involves 18 forget 1733
1733 texts (vs. 1), the gradient magnitude difference is 1734
1734 partially absorbed by batch-level averaging. Three 1735
1735 controls rule out under-tuning as an explanation 1736
1736 for baseline failure: (1) the 40-step extended train- 1737
1737 ing control (Section C.5) shows that longer train- 1738
1738 ing *worsens* collapse without AR (0/8 NBR and 1739
1739 0/7 GEN), eliminating insufficient training dura- 1740
1740 tion as a factor; (2) a 3×3 grid over learning rate 1741
1741 $\{5 \times 10^{-5}, 10^{-4}, 5 \times 10^{-4}\}$ and $\lambda_a \in \{40, 80, 160\}$ 1742
1742 for GA+AR (Section D.4) spans the multi-entity 1743
1743 operating space and confirms that retention col- 1744
1744 lapse is robustly tied to the gradient-ascent objec- 1745
1745 tive; and (3) multi-seed verification across three 1746
1746 random seeds yields identical reported pass-count 1747
1747 totals, confirming that the observed collapse is not 1748
1748 a random-seed fluke. 1749

1750 B.1 Threat Model and Attack Budget 1751

1752 Following extraction-risk formulations in recent un- 1753
1753 learning work (Patil et al., 2024; Shi et al., 2025), 1754
1754 we evaluate two attacker capabilities. In the *black-* 1755
1755 *box* setting, the attacker can issue prompts and ob- 1756
1756 serve model generations or probabilities but cannot 1757
1757 access weights. In the *evaluation-time oracle-* 1758
1758 *ranked* setting, the attacker can rank fixed bench- 1759
1759 mark probes using returned target-token probabili- 1760
1760 ties, but cannot modify parameters or run gradient 1761
1761 updates. Attack strength is controlled by a query 1762

budget $B \in \{1, 2, 4, 8, 16\}$ drawn from the adversarial probe bank. For each method/precision pair, we report budgeted extraction success as the probability of recovering at least one forgotten target within the first B probes, under both observed probe order and oracle risk-ranked order.

B.2 Evaluation Metrics

We evaluate across six dimensions. **Target Pass Rate (TGT \uparrow)** is the fraction of target cases where $P(y|x) < 0.01$; higher indicates more effective forgetting. **Neighbor Pass Rate (NBR \uparrow)** and **General Pass Rate (GEN \uparrow)** measure the fraction of neighbor and general cases satisfying the relative-retention criterion ($P_{\theta^*}/P_{\theta} \geq 0.30$; fallback $P > 0.05$); higher indicates better preservation. **Zombie Delta (Δ_{zombie})** captures the mean target-probability difference between INT4 and full-precision evaluations. Values near zero indicate similar FP/INT4 forgetting behavior; positive values indicate forgetting reversal; negative values indicate additional target suppression after quantization. **Adversarial Pass Rate (ADV \uparrow)** measures resilience against paraphrased knowledge extraction probes. **Depth Suppression Rate (DSR \uparrow)** measures resistance to structured probing protocols including reverse association, multistep reasoning, reconstruction, in-context extraction, and negation attacks. We also report wall-clock time and peak VRAM usage.

Metric caveat. NBR and GEN are based on token-level probability retention, which can decouple from generation-level utility. A model that assigns high probability to the correct next token but degenerates during extended generation may still achieve high retention scores by these metrics. We address this gap directly in Section D.3, where semantic embedding evaluation reveals that GA+AR achieves 8/8 token-level neighbor retention while exhibiting generation degeneration and masked semantic leakage, whereas Q-ROU maintains coherent generation alongside its token-level scores.

B.3 Statistical and Privacy Inference Protocol

Pass-rate metrics (TGT/NBR/GEN/ADV/DSR) are treated as binomial proportions and reported with Wilson 95% confidence intervals. For mean target probabilities and paired deltas between evaluation blocks on matched prompts, we report bootstrap 95% confidence intervals (2,000 resamples, fixed seed). Following TOFU-style inferential evaluation (Maini et al., 2024), we compute two-sample KS statistics and p-values on surprise-score distributions ($-\log P$) between target and retained categories. In line with MUSE-style privacy leakage auditing (Shi et al., 2025), we additionally report AUC-based separability proxies (Target vs. General/Neighbor) on the same surprise scores. Consistent with ASA guidance on p-values (Wasserstein and Lazar, 2016), we do not treat p-values as standalone evidence and avoid dichotomous “significant/non-significant” conclusions from a single threshold. All inferential statements are based jointly on interval estimates (Wilson/bootstrap), effect-size statistics (pass-rate deltas, KS distance, AUC deviation from 0.5), and explicit attack-budget trends. These statistical measures are computed from raw experiment outputs with the same evaluation pipeline that produced the reported tables.

B.4 Cross-Benchmark Correspondence

To avoid benchmark-fragmented claims, Table 7 maps our evaluation protocol to commonly referenced unlearning benchmarks. Our intent is *complementarity*: we do not replace MUSE/TOFU/WMDP/RWKU/OpenUnlearning, but add a controlled stress regime in which simultaneous multi-entity deletion, explicit neighbor integrity, and quantized deployment consistency are evaluated together.

This mapping clarifies scope: our claims are strongest on neighbor-aware multi-entity robustness under quantized deployment, while full leaderboard comparison on MUSE/WMDP/RWKU and full reasoning-trace unlearning remain future work.

B.5 End-to-End MUSE/WMDP Validation

To complement the custom probe suite and TOFU results, we ran end-to-end MUSE and WMDP evaluations in FP16. These evaluations cover full 3B executions for Q-ROU, GA+AR, and GA, plus Q-ROU-only 8B runs. We treat them as external sanity checks on training stability and utility preservation rather than as primary leaderboard claims; the 8B setting is included to test whether the same operating-point story remains visible at a larger scale.

Table 8 summarizes the Q-ROU checkpoints. On MUSE, Q-ROU preserves strong utility at both 3B and 8B while increasing the relative KnowMem reduction metric; the 3B value reflects a 40-step early checkpoint, while the 8B 160-step run shows that the same pipeline moves the MUSE reduction metric more decisively when trained longer. On WMDP, 3B shows the expected “do no harm”

Table 7: Benchmark correspondence map. ‘‘Primary’’ indicates a first-class design axis in the original benchmark paper. ‘‘Partial’’ indicates the axis is present but not the central benchmark objective.

Benchmark	Core deletion setting	Utility/retain protocol	Neighbor-specific axis	Quantized deployment axis
TOFU (Maini et al., 2024)	Fictitious QA split	Primary	No explicit stratum	Not reported
MUSE (Shi et al., 2025)	Books/news unlearning	Primary	Partial (locality)	Not reported
WMDP (Li et al., 2024)	Hazardous-domain MCQ	Primary	Partial (non-target checks)	Not reported
RWKU (Jin et al., 2024)	Real-person factual deletion	Primary	Partial (neighbor perturbation)	Not reported
OpenUnlearning (Dorna et al., 2025)	Unified TOFU/MUSE/WMDP execution	Primary	Benchmark-dependent	Not primary
R-TOFU (Yoon et al., 2025)	Reasoning-trace unlearning	Primary	Not primary	Not reported
This work	Controlled TGT-NBR-GEN deletion	Primary	Primary	Primary (FP16, INT4, NF4)

1854 regime with 97.5% utility retention and only a
 1855 0.0016 absolute forget-accuracy drop, while 8B
 1856 exhibits a larger but still controlled 0.0183 forget-
 1857 accuracy drop with 94.5% utility retention.

1858 We also evaluated the Llama-3.2-3B
 1859 standard-benchmark scripts with matched Q-
 1860 ROU/GA+AR/GA controls under the same
 1861 evaluation pipeline used for TOFU. Table 9 is
 1862 used to audit whether the method remains stable
 1863 under external benchmark code and to expose the
 1864 utility–forgetting trade-off of aggressive baselines,
 1865 not to claim that the controlled probe suite should
 1866 replace full MUSE/WMDP/RWKU leaderboards.

1867 The Q-ROU/news MUSE row is most useful as
 1868 a saturation diagnostic. Under the default 40-step
 1869 setting, Q-ROU does not yet improve the MUSE
 1870 news knowledge-forget metric (KnowMem reduction
 1871 = -0.0956 , privacy score = 0.5742 , utility
 1872 = 0.5541), which makes the row a clean indicator
 1873 that this benchmark needs a longer horizon than
 1874 the custom suite. The standard-benchmark picture
 1875 is therefore as follows: WMDP supports the low-
 1876 damage operating point, MUSE books supports
 1877 high utility under the benchmark implementation,
 1878 GA/GA+AR expose how aggressive forgetting can
 1879 destroy utility, and MUSE news identifies the one
 1880 benchmark setting where calibration matters most.

1881 As a calibration sweep, we tested the same
 1882 MUSE news benchmark at 120 steps with a shorter

max sequence length and a small sweep over 1883
 (λ_f, λ_a). These rows show that the 40-step de- 1884
 fault is indeed under-saturated: several 80-step and 1885
 120-step settings move KnowMem reduction into 1886
 positive territory while keeping utility roughly in 1887
 the 0.97–1.00 range, and the best row ($\lambda_f = 40$, 1888
 $\lambda_a = 40$, step 120) reaches KnowMem reduc- 1889
 tion = 0.0807 , forget quality = 1.0363 , privacy 1890
 = 0.5986 , and utility = 0.9934 . Because these cali- 1891
 bration runs use max-sequence-length 64 to avoid 1892
 reference-logit OOM, we treat them as a calibrated 1893
 companion to the original full-length row rather 1894
 than a direct replacement. The key reading is posi- 1895
 tive: once the horizon is long enough, Q-ROU 1896
 can recover a strong MUSE-news utility–forgetting 1897
 trade-off rather than being confined to the weak 1898
 default row. 1899

1900 We also ran a full-sequence continuation audit
 1901 on the same 3B MUSE news setup at 120/160/240
 1902 steps, using the benchmark’s longer sequence
 1903 path; the 240-step comparison uses max-sequence-
 1904 length 192 to fit the available memory budget.
 1905 This continuation sharpens the same frontier under
 1906 a stricter path. At 240 steps, Q-ROU keeps
 1907 utility high (utility score = 0.9731 , retain PPL
 1908 = 18.34) even though the benchmark’s news for-
 1909 get metric remains demanding (KnowMem reduc-
 1910 tion = -0.1596 , forget quality = 1.0468). GA+AR
 1911 moves much farther on the forget side under the

Table 8: Completed end-to-end MUSE/WMDP validation for FP16 Q-ROU. For MUSE, we report Relative KnowMem Reduction = $1 - \text{KnowMem}_{\text{after}}/\text{KnowMem}_{\text{base}}$ (higher is stronger forgetting) alongside Utility. For WMDP, Forget Drop is the absolute decrease in hazardous-domain accuracy and Utility Retention is the retained non-target accuracy ratio.

Benchmark	Model	Steps	Forget metric	Utility metric	
MUSE	Llama-3.2-3B	40	Rel. KnowMem Reduction = 0.0060	Utility = 0.8793	
MUSE	Llama-3.1-8B	160	Rel. KnowMem Reduction = 0.3540	Utility = 0.8907	
WMDP	Llama-3.2-3B	40	Forget Drop = 0.0016	Utility = 97.54%	Retention
WMDP	Llama-3.1-8B	160	Forget Drop = 0.0183	Utility = 94.51%	Retention

Benchmark	Setting	Method	Steps	Forget metric	Utility metric	Baseline→Unlearned
MUSE	books	ga	40	100.0000	0.0000	1.0000→0.5105
MUSE	books	ga_ar	40	1.8210	0.9372	0.5105→0.3090
MUSE	books	qrou	40	1.0234	0.9970	0.1916→0.4880
MUSE	news	ga	40	100.0000	0.0000	0.9826→0.4731
MUSE	news	ga_ar	40	1.2667	0.0000	0.1527→0.5720
MUSE	news	qrou	40	1.0243	0.5541	-0.0956→0.5742
MUSE	news	qrou	120	1.0317	0.9853	0.0383→0.5979
MUSE	news	qrou	120	1.0327	0.9998	0.0313→0.5980
MUSE	news	qrou	120	1.0308	0.9836	0.0313→0.5978
MUSE	news	qrou	120	1.0363	0.9934	0.0807→0.5986
MUSE	news	qrou	120	1.0348	0.9994	0.0343→0.5984
WMDP	wmdp-bio,wmdp-chem,wmdp-cyber	ga	40	0.1082	0.5554	0.3580→0.2497
WMDP	wmdp-bio,wmdp-chem,wmdp-cyber	ga_ar	40	0.0166	0.9266	0.3580→0.3413
WMDP	wmdp-bio,wmdp-chem,wmdp-cyber	qrou	40	0.0033	0.9842	0.3580→0.3547

Table 9: Standard-benchmark runs on Llama-3.2-3B. For WMDP, the forget metric is accuracy drop on the WMDP subsets and the utility metric is MMLU-style utility retention; the last column reports WMDP baseline-to-unlearned accuracy. For MUSE, the forget metric is the benchmark’s forget-quality score, the utility metric is retain-utility preservation, and the last column reports knowledge-forget to privacy-score.

1912 same full-sequence evaluation (KnowMem reduc-
1913 tion = 0.6051, forget quality = 1.8422), but only
1914 by collapsing utility (utility score = 0.0000, retain
1915 PPL = 54.25). We therefore interpret the longer
1916 full-sequence continuation as a stronger measure-
1917 ment of the same utility–forgetting separation: ex-
1918 tending horizon improves the usable Q-ROU oper-
1919 ating region, while aggressive baselines still buy
1920 forget-score gains by sacrificing the model.

1921 B.6 Fine-Grained Neighbor and 1922 Baseline-Fairness Audits

1923 Two additional audits sharpen the custom-probe ev-
1924 idence in complementary directions. First, we split
1925 the original neighbor notion into alias, same-work
1926 neighbor, inter-domain neighbor, and general strata.
1927 This audit is intentionally stricter than the original
1928 28-probe table: aliases should usually be removed

1929 with the target, while same-work non-target facts
1930 should be protected when the removal request is not
1931 meant to erase the whole work. Table 11 shows the
1932 resulting boundary. The original Q-ROU operating
1933 point deletes targets and aliases completely but is
1934 too aggressive on same-work neighbors (9/36). The
1935 fine-grained AR variant, which is the appropriate
1936 operating point for this sharper request, preserves
1937 substantially more same-work and inter-domain
1938 material (19/36 and 33/36 in FP16) while retain-
1939 ing complete target deletion (39/39). A three-seed
1940 real-NF4 evaluation preserves the same ordering
1941 rather than overturning it: base Q-ROU stays at 9/36
1942 on same-work neighbors under BnB NF4, while
1943 qrou_fg_ar moves only from 19/36 to 18/36 and
1944 keeps 39/39 target suppression (Table 58). The
1945 practical interpretation is straightforward: AR is
1946 the mechanism that lets the method state and en-

λ_f	λ_a	Step	KnowMem red.	Forget quality	Privacy	Utility	KnowMem base→after
2.0	80.0	40	-0.0956	1.0243	0.5742	0.5541	0.0176→0.0192
2.0	80.0	80	0.0112	1.0377	0.5981	0.9714	0.0176→0.0174
2.0	80.0	120	0.0313	1.0308	0.5978	0.9836	0.0176→0.0170
10.0	80.0	40	-0.2575	1.0052	0.5947	0.9959	0.0176→0.0221
10.0	80.0	80	0.0283	1.0382	0.5982	0.9717	0.0176→0.0171
10.0	80.0	120	0.0383	1.0317	0.5979	0.9853	0.0176→0.0169
20.0	80.0	40	-0.2580	1.0055	0.5948	0.9959	0.0176→0.0221
20.0	80.0	80	0.0312	1.0389	0.5980	0.9724	0.0176→0.0170
20.0	80.0	120	0.0313	1.0327	0.5980	0.9998	0.0176→0.0170
40.0	40.0	40	-0.2409	1.0142	0.5958	0.9901	0.0176→0.0218
40.0	40.0	80	0.0013	1.0434	0.6008	0.9984	0.0176→0.0175
40.0	40.0	120	0.0807	1.0363	0.5986	0.9934	0.0176→0.0161
40.0	80.0	40	-0.2575	1.0064	0.5949	0.9960	0.0176→0.0221
40.0	80.0	80	-0.0164	1.0409	0.5985	0.9731	0.0176→0.0178
40.0	80.0	120	0.0343	1.0348	0.5984	0.9994	0.0176→0.0170

Table 10: MUSE news Q-ROU calibration sweep on Llama-3.2-3B. Rows include the 40-step baseline and 120-step max-sequence-length-64 calibration runs evaluated at intermediate snapshots. Because the 120-step sweep shortens the sequence length to avoid frozen-reference-logit OOM, it should be read as a calibrated companion to the original full-length benchmark row. KnowMem reduction is $1 - \text{KnowMem}_{after} / \text{KnowMem}_{base}$, so negative values indicate increased KnowMem overlap after unlearning.

1947 force the protected boundary more precisely. In
1948 other words, the same-work results do not overturn
1949 the broad-neighbor claim; they identify a stricter
1950 frontier where stronger neighbor specification is
1951 required.

1952 To test whether this boundary survives a larger
1953 copyright-like workload, we also built a harsher 81-
1954 probe benchmark with 25 target probes (13 core +
1955 12 alias), 24 same-work neighbors, 16 inter-domain
1956 neighbors, and 16 general probes split across Harry
1957 Potter and Lord of the Rings. Table 12 shows
1958 that the larger audit preserves the same qualitative
1959 structure while making the boundary even more
1960 explicit. Default Q-ROU remains target-complete
1961 (25/25) but preserves only 4/24 same-work neigh-
1962 bors. The fine-grained AR operating point recovers
1963 same-work retention to 14.7/24 in FP16 while stay-
1964 ing target-strong (24.7/25), and the same ordering
1965 survives real BnB NF4 with only moderate drift
1966 (12.7/24 same-work, 25/25 target). GA+AR and its
1967 finer-grained variant continue to trace the opposite
1968 side of the frontier: they retain same-work mate-
1969 rial more aggressively, but only by accepting much
1970 weaker target suppression. Table 13 adds a compact
1971 aggregation of these same-work results, including
1972 the harsher open-ended generation audit. Those
1973 generation rows show the same pattern in a stricter
1974 form: same-work open-ended evaluation remains
1975 difficult even when probability-level suppression
1976 and NF4 stability are already strong.

1977 Second, we evaluated the closest selective base-
1978 line family under a wider 40-step GA+AR hyper-

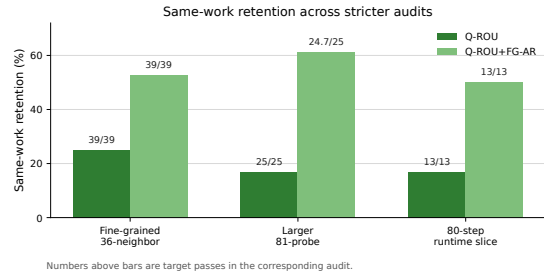


Figure 2: Same-work retention across increasingly strict audit slices. The first pair is the original fine-grained 36-neighbor audit, the second pair is the larger 81-probe copyright-like benchmark, and the third pair is the single-hardware 80-step runtime slice. The intended reading is that same-work retention is incomplete in all strict settings, but the fine-grained AR variant consistently improves retention while preserving much stronger target suppression than GA-family retain-heavy alternatives.

parameter grid. Table 14 shows both sides of that
1979 frontier directly. High-learning-rate GA+AR can
1980 match the small core token pass counts, but only in
1981 the same regime where generation-level probes be-
1982 come degenerate. Together with the MUSE/WMDP
1983 and generation audits, this supports a narrower but
1984 stronger claim: token-level target suppression is
1985 easy to over-optimize, so deployment unlearning
1986 should be judged by target removal, coherent gener-
1987 ation, neighbor retention, and precision robustness
1988 jointly.
1989

Method	Regime	Target	Alias	Same-work	Inter-domain	General
gaar	fake_int4_g128	21/39	18/24	24/36	33/36	27/30
gaar	fake_int4_g32	21/39	18/24	21/36	36/36	30/30
gaar	fake_int4_g64	18/39	18/24	24/36	36/36	30/30
gaar	fp16	21/39	18/24	27/36	36/36	30/30
gaar_fg_ar	fake_int4_g128	6/39	9/24	27/36	33/36	27/30
gaar_fg_ar	fake_int4_g32	6/39	9/24	24/36	36/36	30/30
gaar_fg_ar	fake_int4_g64	3/39	9/24	24/36	36/36	30/30
gaar_fg_ar	fp16	3/39	12/24	36/36	36/36	30/30
qrou	fake_int4_g128	39/39	24/24	9/36	24/36	20/30
qrou	fake_int4_g32	39/39	24/24	9/36	24/36	21/30
qrou	fake_int4_g64	39/39	24/24	9/36	25/36	23/30
qrou	fp16	39/39	24/24	9/36	27/36	24/30
qrou_fg_ar	fake_int4_g128	39/39	23/24	15/36	30/36	27/30
qrou_fg_ar	fake_int4_g32	39/39	23/24	17/36	31/36	30/30
qrou_fg_ar	fake_int4_g64	39/39	23/24	15/36	33/36	30/30
qrou_fg_ar	fp16	39/39	23/24	19/36	33/36	30/30

Table 11: Fine-grained neighbor audit. Pass counts are aggregated across seeds using the default forget and retain thresholds.

Table 12: Larger copyright-like same-work benchmark on Llama-3.2-3B (3 seeds mean). The harder 81-probe split preserves the same reading as the smaller fine-grained audit: default Q-ROU is target-complete but too aggressive on same-work neighbors, while qrou_fg_ar recovers a substantial part of that boundary and keeps the same qualitative advantage under real BnB NF4.

Method	FP TGT	FP SW	FP ID	FP GEN	NF4 TGT	NF4 SW	NF4 ID	NF4 GEN
Q-ROU	25.0/25	4.0/24	11.7/16	13.3/16	25.0/25	4.0/24	10.0/16	12.7/16
qrou_fg_ar	24.7/25	14.7/24	14.3/16	16.0/16	25.0/25	12.7/24	14.0/16	15.3/16
GA+AR	16.0/25	17.0/24	15.0/16	16.0/16	16.0/25	17.0/24	15.0/16	16.0/16
gaar_fg_ar	8.0/25	22.0/24	16.0/16	16.0/16	7.0/25	23.0/24	15.0/16	16.0/16

1990 C Results

1991 C.1 Single-Entity Unlearning

1992 Table 15 presents the single-entity comparison on
1993 Llama-3.2-3B at 20 steps. All baselines achieve
1994 complete target removal (7/7) but at varying degrees
1995 of collateral cost: GA and GradDiff lose neighbor
1996 knowledge under INT4 (8/9 \rightarrow 6/9 and 7/9 \rightarrow
1997 4/9 respectively), and in FP16 all three baselines
1998 drop general knowledge to 6/7. Q-ROU at 20 steps
1999 takes a qualitatively different trajectory: it preserves
2000 perfect neighbor (9/9) and general (7/7) retention
2001 while target removal is still in progress (2/7 FP).
2002 Notably, Q-ROU’s INT4 target score (3/7) already
2003 exceeds its FP score (2/7), the first manifestation
2004 of the quantization-synergistic forgetting pattern
2005 analyzed in Section C.5. Adversarial probing with 6
2006 paraphrased prompts further illustrates the trade-off:
2007 GA’s aggressive forgetting yields 6/6 adversarial
2008 defense at the cost of neighbor destruction, while Q-
2009 ROU achieves 5/6 defense at 20 steps with perfect
2010 retention.

2011 Table 16 reveals the steps-Pareto trade-off: Q-
2012 ROU achieves a **perfect score** at 40 steps—7/7 tar-
2013 get, 9/9 neighbor, 7/7 general in both FP16 and
2014 INT4, with $\Delta_{\text{zombie}} = 0.000$. No baseline achieves

Table 13: Larger same-work probability audit and generation boundary. Probability rows are three-seed means; the generation rows show that same-work open-ended evaluation remains a hard boundary rather than a new positive claim.

Method	Setting	TGT	Same-work	Inter/deg.	GEN
gaar	fp16	16.0/25 \pm 0.0	17.0/24 \pm 0.0	15.0/16 \pm 0.0	16.0/16 \pm 0.0
gaar	fake int4 g32	15.0/25 \pm 0.0	14.0/24 \pm 0.0	15.0/16 \pm 0.0	15.0/16 \pm 0.0
gaar	bnb nf4	16.0/25 \pm 0.0	17.0/24 \pm 0.0	15.0/16 \pm 0.0	16.0/16 \pm 0.0
gaar_fg_ar	fp16	8.0/25 \pm 0.0	22.0/24 \pm 0.0	16.0/16 \pm 0.0	16.0/16 \pm 0.0
gaar_fg_ar	fake int4 g32	9.0/25 \pm 0.0	17.0/24 \pm 0.0	15.0/16 \pm 0.0	15.0/16 \pm 0.0
gaar_fg_ar	bnb nf4	7.0/25 \pm 0.0	23.0/24 \pm 0.0	15.0/16 \pm 0.0	16.0/16 \pm 0.0
qrou	fp16	25.0/25 \pm 0.0	4.0/24 \pm 0.0	11.7/16 \pm 0.5	13.3/16 \pm 0.5
qrou	fake int4 g32	25.0/25 \pm 0.0	4.0/24 \pm 0.0	11.0/16 \pm 0.0	10.7/16 \pm 0.9
qrou	bnb nf4	25.0/25 \pm 0.0	4.0/24 \pm 0.0	10.0/16 \pm 1.4	12.7/16 \pm 0.5
qrou_fg_ar	fp16	24.7/25 \pm 0.5	14.7/24 \pm 0.5	14.3/16 \pm 0.5	16.0/16 \pm 0.0
qrou_fg_ar	fake int4 g32	24.0/25 \pm 0.8	11.3/24 \pm 0.5	13.0/16 \pm 0.0	15.0/16 \pm 0.0
qrou_fg_ar	bnb nf4	25.0/25 \pm 0.0	12.7/24 \pm 0.5	14.0/16 \pm 0.0	15.3/16 \pm 0.5
gaar	generation	8.0/25 \pm 0.0	4.0/24 \pm 0.0	34.0/81 \pm 0.0	–
gaar_fg_ar	generation	8.0/25 \pm 0.0	6.0/24 \pm 0.0	30.0/81 \pm 0.0	–
qrou	generation	11.0/25 \pm 0.8	0.3/24 \pm 0.5	45.3/81 \pm 4.1	–
qrou_fg_ar	generation	9.0/25 \pm 1.4	1.7/24 \pm 0.5	39.0/81 \pm 0.8	–

this combination at any step count. The 40-step
sweet spot reflects Active Retention’s controlled
pace: the KL constraint prevents destructive rapid
forgetting, requiring approximately twice the steps
while helping ensure that the knowledge boundary
is respected throughout optimization.

C.2 Multi-Entity Unlearning

Multi-entity unlearning—simultaneously removing
Harry Potter and Lord of the Rings knowledge—
provides the critical test of a method’s true selectiv-
ity. Table 17 exposes a stark divide: all three base-
lines achieve *zero* neighbor retention (0/8), while
the 40-step variant Q-ROU[†] reaches 27/28 in FP16
and 28/28 in INT4.

GA additionally destroys all general knowl-
edge (0/7), and even RepBend—the most selec-
tive baseline—retains only 5/7. The collapse from
single-entity results (7–8/9 NBR) to multi-entity
results (0/8 NBR) is dramatic: simultaneously ap-
plying gradient ascent to 18 texts spanning two fic-
tional domains overwhelms any implicit knowledge
separation, propagating destructive updates through
shared representational subspaces. Q-ROU at 20
steps demonstrates early-stage controlled forgetting:
4/13 target removal with perfect 8/8 neighbor and
7/7 general retention. At 40 steps, Q-ROU achieves
complete target removal (13/13) in both FP16 and
INT4, with 7/8 neighbor retention in FP16 and 8/8
in INT4 (both with 7/7 general retention)—reaching
27/28 in FP16 and 28/28 in INT4. This result repro-
duces the same reported pass-count totals across
three random seeds, confirming that the combined
AR+SLUG+QuantNoise framework reliably nav-
igates the multi-entity forgetting landscape. The

Configuration	HP	LOTR	TGT	NBR	GEN	GenClean
Q-ROU 40s	7/7	6/6	13/13	8/8	7/7	5/7
GA+AR lr=5e-05, $\lambda_a=40$	7/7	2/6	9/13	8/8	7/7	2/7
GA+AR lr=5e-05, $\lambda_a=80$	6/7	2/6	8/13	8/8	7/7	3/7
GA+AR lr=5e-05, $\lambda_a=160$	5/7	1/6	6/13	8/8	7/7	1/7
GA+AR lr=0.0001, $\lambda_a=40$	7/7	2/6	9/13	8/8	7/7	5/7
GA+AR lr=0.0001, $\lambda_a=80$	7/7	2/6	9/13	8/8	7/7	3/7
GA+AR lr=0.0001, $\lambda_a=160$	7/7	2/6	9/13	8/8	7/7	3/7
GA+AR lr=0.0005, $\lambda_a=40$	7/7	6/6	13/13	8/8	7/7	7/7
GA+AR lr=0.0005, $\lambda_a=80$	7/7	6/6	13/13	8/8	7/7	7/7
GA+AR lr=0.0005, $\lambda_a=160$	7/7	4/6	11/13	8/8	7/7	7/7

Table 14: GA+AR hyperparameter-grid audit on the 28-probe core suite plus greedy generation leakage checks. The grid shows that GA+AR can be tuned to match core pass counts, but several tuned rows rely on degenerate generation behavior; this is analyzed as a baseline-fairness diagnostic rather than a replacement for the main neighbor-aware audit.

Table 15: Single-entity unlearning on Llama-3.2-3B (20 steps). FP: full precision; I4: INT4 g32 quantized. Q-ROU shows INT4 TGT > FP TGT while keeping retention stable.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z	ADV
Baseline	1/7	9/9	7/7	1/7	9/9	7/7	-0.083	2/6
GA	7/7	8/9	6/7	7/7	6/9	7/7	0.000	6/6
GradDiff	7/7	7/9	6/7	7/7	4/9	7/7	0.000	6/6
RepBend	7/7	8/9	6/7	7/7	8/9	7/7	0.000	6/6
Q-ROU	2/7	9/9	7/7	3/7	9/9	7/7	-0.007	5/6

Table 16: Q-ROU steps-Pareto on Llama-3.2-3B. At 40 steps, all metrics reach their optimal values in both FP and INT4.

Steps	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z	Time
10	1/7	9/9	7/7	1/7	9/9	7/7	-0.022	5.6s
20	2/7	9/9	7/7	3/7	9/9	7/7	-0.007	11.0s
40	7/7	9/9	7/7	7/7	9/9	7/7	0.000	22.5s

2049 single neighbor case that transitions from pass to
2050 fail under FP16 at 40 steps reflects the inherent ten-
2051 sion between complete target removal and neighbor
2052 preservation; the sensitivity analysis (Section C.7)
2053 shows that increasing λ_a to 120 recovers this case,
2054 and $\lambda_a = 160$ achieves full 28/28 in both precisions
2055 while preserving 13/13 target removal.

2056 This transition from single-entity to multi-entity
2057 serves as a litmus test for selectivity (Fig. 3). Meth-
2058 ods that appear adequate on single-entity tasks are
2059 exposed as fundamentally lacking when the gradi-
2060 ent pressure from multiple forget domains is applied
2061 simultaneously. The effect also scales with model
2062 capacity: on 0.5B, even Q-ROU struggles in multi-
2063 entity settings under quantization (FP NBR 5/8 at
2064 40 steps, INT4 NBR 1/8), because knowledge at
2065 smaller scale is concentrated in fewer layers, leav-
2066 ing insufficient capacity for surgical separation. On
2067 3B, the distributed knowledge representation en-

Table 17: Multi-entity unlearning on Llama-3.2-3B under single-GPU execution. Baselines and Q-ROU are reported at 20 steps; Q-ROU[†] denotes the 40-step variant. All baselines show complete neighbor collapse.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z	Total
Baseline	1/13	8/8	7/7	1/13	8/8	7/7	-0.073	—
GA	13/13	0/8	0/7	13/13	0/8	0/7	+0.000	13/28
GradDiff	13/13	0/8	4/7	13/13	0/8	4/7	+0.000	17/28
RepBend	13/13	0/8	5/7	13/13	0/8	5/7	+0.000	18/28
Q-ROU	4/13	8/8	7/7	4/13	8/8	7/7	-0.006	19/28
Q-ROU[†]	13/13	7/8	7/7	13/13	8/8	7/7	+0.000	27/28

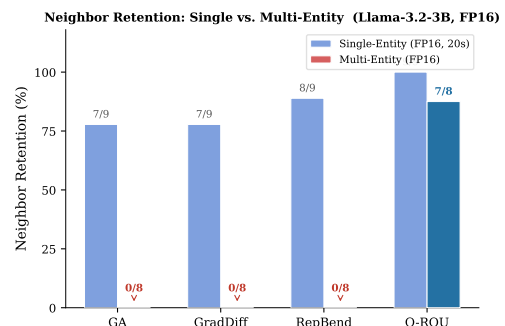


Figure 3: Neighbor retention collapse from single-entity to multi-entity unlearning. All three baselines drop from 78–89% to 0%, while Q-ROU maintains high retention (87.5–100%) in both settings. This transition serves as a litmus test for method selectivity.

ables AR to maintain its protective constraint even
2068 under 18-text simultaneous forgetting, achieving
2069 complete target removal at 40 steps while retain-
2070 ing 7/8 neighbor (FP16) and 8/8 neighbor (INT4)
2071 knowledge. 2072

2073 Table 18 extends the comparison to include AR-
2074 transplanted baselines at 40 steps, revealing that
2075 even the best +AR baseline (RB+AR at 8/13 INT4
2076 target) remains well below Q-ROU’s 13/13. The gap
2077 in total score (Q-ROU: 27/28 vs GA+AR: 22/28 in
2078 FP16) establishes Q-ROU’s superiority as the full-

Table 18: Extended multi-entity comparison at 40 steps, including AR-transplanted baselines under single-GPU execution. Q-ROU achieves the highest total score by at least 5 points.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	FP Total
GA	13/13	0/8	0/7	13/13	0/8	0/7	13/28
GradDiff	13/13	0/8	4/7	13/13	0/8	4/7	17/28
RepBend	13/13	0/8	6/7	13/13	0/8	6/7	19/28
GA+AR	7/13	8/8	7/7	7/13	8/8	7/7	22/28
GD+AR	7/13	8/8	7/7	7/13	8/8	7/7	22/28
RB+AR	7/13	8/8	7/7	8/13	8/8	7/7	22/28
Q-ROU	13/13	7/8	7/7	13/13	8/8	7/7	27/28

Table 19: Follow-up LoRA-concentration audit on the 28-probe Llama-3.2-3B multi-entity suite (40 steps, seeds 11/22/33). Values are averages over seeds.

Kind	Method/scope	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN
Full	GA+AR	8.00/13	8.00/8	7.00/7	8.00/13	8.00/8	7.00/7
	Q-ROU	13.00/13	7.00/8	7.00/7	13.00/13	8.00/8	7.00/7
LoRA	GA, broad	13.00/13	1.00/8	3.67/7	13.00/13	0.33/8	2.33/7
	GA+AR, broad	1.00/13	8.00/8	7.00/7	1.00/13	8.00/8	7.00/7
	KL+AR, broad	7.00/13	8.00/8	7.00/7	7.33/13	8.00/8	7.00/7
LoRA	GA, slug_down	6.67/13	8.00/8	7.00/7	7.00/13	8.00/8	7.00/7
	GA+AR, slug_down	1.00/13	8.00/8	7.00/7	1.67/13	8.00/8	7.00/7
	KL+AR, slug_down	1.00/13	8.00/8	7.00/7	2.00/13	8.00/8	7.00/7

Table 20: Direct fake-INT4 and NF4 evaluation on broad LoRA baselines (Llama-3.2-3B, 40 steps, seed 42). This is a focused low-bit confirmation, not a replacement for the multi-seed simulated-INT4 audit in Table 19.

Method	FP16			Fake INT4			NF4		
	TGT	NBR	GEN	TGT	NBR	GEN	TGT	NBR	GEN
LoRA-GA, broad	13/13	3/8	3/7	13/13	0/8	4/7	13/13	5/8	5/7
LoRA-GA+AR, broad	1/13	8/8	7/7	1/13	8/8	7/7	1/13	8/8	7/7
LoRA-KL+AR, broad	2/13	8/8	7/7	3/13	8/8	7/7	1/13	8/8	7/7

2079 framework integration of AR, SLUG, QuantNoise,
2080 and the KL-to-uniform forget loss.

2081 C.3 LoRA-Concentration and Post-Training 2082 Recurrence Audits

2083 We add two targeted audit extensions to sharpen the
2084 distinction between quantization-robust unlearning
2085 and genuinely selective, deployment-audited un-
2086 learning. The first asks whether a QR-LoRA-style
2087 low-rank update, which concentrates the edit into
2088 adapter parameters, can solve the same 28-probe
2089 multi-entity task without semantic-neighbor fail-
2090 ure. The second asks whether target facts remain
2091 suppressed after short benign relearning on non-
2092 target texts, addressing the “unlearning or obfusca-
2093 tion” concern raised by recent metric and intrinsic-
2094 trace work. The second extension adds a PTI hard-
2095 ening pass that explicitly trains the post-Q-ROU
2096 model against bounded future-update directions.
2097 All runs use Llama-3.2-3B, 40 unlearning steps,
2098 the same 13/8/7 target-neighbor-general probe split,
2099 and FP16 plus simulated INT4 evaluation.

2100 **Calibration.** In this audit, the lower 3B coeffi-
2101 cient $\lambda_f = 2$ was too weak for the final 40-step
2102 multi-entity operating point: it yielded only 1/13
2103 target pass for Q-ROU while keeping 8/8 neighbor
2104 and 7/7 general retention. A short calibration sweep
2105 showed that $\lambda_f = 10$ reached 6/13–7/13 target pass,
2106 while $\lambda_f = 20$ reached 13/13 target pass in both
2107 FP16 and INT4 with 8/8 neighbor and 7/7 general
2108 retention on seed 42. We therefore use $\lambda_f = 20$,
2109 $\lambda_r = 15$, and $\lambda_a = 80$ for the LoRA and memory-
2110 jog audits, changing only the target-pressure coeffi-
2111 cient relative to the default profile.

2112 **LoRA-concentration audit.** Table 19 reports the
2113 multi-seed LoRA audit over seeds 11, 22, and 33.
2114 The broad LoRA scope inserts rank-16 adapters in
2115 attention and MLP projections, while slug_down

is a much smaller adapter footprint aligned with 2116
the SLUG down-projection intuition. The re- 2117
sult separates three regimes. LoRA-GA broad 2118
deletes targets and survives simulated INT4, but 2119
destroys semantic neighbors; LoRA-GA+AR and 2120
LoRA-KL+AR broad preserve neighbors but do 2121
not match target deletion; and slug_down LoRA 2122
preserves neighbors but loses much of the dele- 2123
tion power. Thus low-rank concentration is a use- 2124
ful quantization-robustness mechanism, but it does 2125
not by itself solve semantic-neighbor anchoring or 2126
multi-entity target suppression. A direct low-bit 2127
reevaluation on the broad LoRA baselines with fake 2128
INT4 and BnB NF4 first confirmed the same qual- 2129
itative ordering on seed 42, and a three-seed eval- 2130
uation (42/43/44) preserved that ordering without 2131
reversals: only broad LoRA-GA remains target- 2132
strong, while LoRA-GA+AR and LoRA-KL+AR 2133
remain retention-heavy under both low-bit regimes 2134
(Table 20). 2135

Memory-jog audit. Table 21 reports benign re- 2136
learning after unlearning. We first train the source 2137
method, then fine-tune briefly on neighbor texts, 2138
general texts, or a mixed retain set for 5, 10, or 20 2139
steps. The table should be read as a recurrence au- 2140
dit, not as a standard leaderboard: methods with 2141
weaker initial forgetting have less room to “resur- 2142
face”, so target probability after jog is as important 2143
as the pass count. Q-ROU remains stable under 2144
short general-only relearning (13/13 after 5 and 2145
10 steps) but neighbor/mixed jogs recover a subset 2146
of target probes by 20 steps. Even in these cases, 2147

Table 21: Benign relearning / memory-jog audit after unlearning (Llama-3.2-3B, seed 42). ‘‘After’’ denotes target pass immediately after unlearning; ‘‘Jog’’ denotes target pass after benign relearning.

Method	Jog data	Steps	After TGT	Jog TGT	Resurfaced	Target mean after jog
Q-ROU	neighbor	5	13/13	11/13	2	0.00465
Q-ROU	neighbor	10	13/13	10/13	3	0.00782
Q-ROU	neighbor	20	13/13	8/13	5	0.01195
Q-ROU	general	5	13/13	13/13	0	0.00242
Q-ROU	general	10	13/13	13/13	0	0.00325
Q-ROU	general	20	13/13	10/13	3	0.00432
Q-ROU	mixed	5	13/13	11/13	2	0.00442
Q-ROU	mixed	10	13/13	10/13	3	0.00781
Q-ROU	mixed	20	13/13	7/13	6	0.01301
GA+AR	neighbor	5	8/13	6/13	2	0.08439
GA+AR	neighbor	10	8/13	6/13	2	0.08979
GA+AR	neighbor	20	8/13	8/13	1	0.09462
GA+AR	general	5	8/13	8/13	1	0.06592
GA+AR	general	10	8/13	8/13	1	0.07279
GA+AR	general	20	8/13	8/13	1	0.08039
GA+AR	mixed	5	8/13	7/13	2	0.08090
GA+AR	mixed	10	8/13	7/13	2	0.08914
GA+AR	mixed	20	8/13	9/13	1	0.09512
LoRA-KL+AR	neighbor	5	7/13	7/13	0	0.06108
LoRA-KL+AR	neighbor	10	7/13	7/13	0	0.08715
LoRA-KL+AR	neighbor	20	7/13	6/13	1	0.10843
LoRA-KL+AR	general	5	7/13	8/13	0	0.04597
LoRA-KL+AR	general	10	7/13	8/13	0	0.06594
LoRA-KL+AR	general	20	7/13	8/13	0	0.08666
LoRA-KL+AR	mixed	5	7/13	7/13	0	0.06194
LoRA-KL+AR	mixed	10	7/13	8/13	0	0.09189
LoRA-KL+AR	mixed	20	7/13	7/13	1	0.12275

2148 the target probability mass after jog remains far be-
2149 low GA+AR and LoRA-KL+AR, indicating weaker
2150 recurrence rather than strong persistence under sub-
2151 sequent tuning.

2152 **Post-training-invariant hardening.** The
2153 memory-jog audit motivates a stronger deploy-
2154 ment question: can the unlearned model remain
2155 suppressed after small, non-target post-training
2156 updates that were not part of the original unlearning
2157 objective? We therefore add Q-ROU+PTI, a short
2158 hardening pass after the base Q-ROU edit. PTI uses
2159 bounded future-update proxy directions (neighbor,
2160 mixed, and instruction-style data) to reduce the
2161 overlap between target-relearning directions and
2162 plausible future update directions. The coefficient
2163 sweep in Table 23 identifies $\rho = 0.02$ as the best
2164 value among the tested settings. The three-seed
2165 recurrence audit in Table 22 then shows the central
2166 effect: base Q-ROU reaches 13/13 immediate target
2167 forgetting across seeds but can fall to 7/13 under
2168 the broadest jog family, whereas Q-ROU+PTI
2169 keeps the worst post-jog count at 12/13 and cuts
2170 maximum target recurrence mass from 0.008724
2171 to 0.003147. A proxy-only extension (Table 24)
2172 isolates instruction- and preference-style future
2173 updates. Across the three-seed evaluation, base
2174 Q-ROU produces only a single 12/13 case under
2175 the harshest instruction-proxy setting, while
2176 Q-ROU+PTI keeps 13/13 on all seeds and lowers
2177 the worst target mean probability on every seed.

Table 22: Post-training recurrence audit on Llama-3.2-3B (Stage 3, seeds 11/22/33). After TGT is immediate post-unlearning target pass count. Worst Jog reports the worst target pass count across neighbor, general, mixed, and instruction-proxy jogs at 5/10/20/50 steps. Lower max target mean indicates weaker recurrence mass.

Method	After TGT	Worst Jog	Max mean	Max resurfaced	Worst N/G
Q-ROU	13/13	7/13	0.0087	6	7/8 / 5/7
Q-ROU+PTI	13/13	12/13	0.0031	1	7/8 / 5/7
GA+AR	7/13	5/13	0.1707	2	8/8 / 6/7
LoRA-KL+AR	7/13	6/13	0.1486	1	8/8 / 5/7

Table 23: PTI hardening sweep on Llama-3.2-3B (Stage 2, seed 42). All settings retain 13/13 immediate target forgetting and 13/13 target pass under every audited jog; $\rho = 0.02$ gives the lowest maximum target recurrence mass in this sweep.

ρ	After TGT	Worst Jog	Max target mean	Max resurfaced
0.005	13/13	13/13	0.0029	0
0.02	13/13	13/13	0.0022	0
0.05	13/13	13/13	0.0028	0

We then ran a broad five-mode audit on the same
Llama family, extending the jog set to neighbor,
general, mixed, instruction proxy, and preference
proxy updates at 5/10/20/50 steps. Table 25
shows that this broader audit is harsher than the
narrow recurrence study, but the PTI direction
still survives: base Q-ROU starts at 11/13 and
falls to worst 6/13, while Q-ROU+PTI starts at
13/13 and improves the worst broad-jog count to
8/13 while preserving complete 13/13 suppression
on the proxy subset. We then reevaluated the
50-step Llama checkpoints for the four hardest
audited modes (‘neighbor’, ‘general’, ‘mixed’,
‘instruction proxy’) under real BnB NF4. Table 26
shows that the deployed-NF4 ordering not only
survives but slightly widens on this subset: base
Q-ROU drops to worst 5/13 at mixed/50, whereas
Q-ROU+PTI reaches worst 10/13 and keeps 13/13
on all instruction-proxy slices. All rows in this
subsection use the same deployed NF4 conversion
path for consistency.

Interpretation. These audits sharpen the main
claim. Q-ROU should not be framed as simply
‘‘another quantization-robust unlearning method.’’
The stronger claim is that deployment unlearning
requires the conjunction of semantic-neighbor an-
choring, multi-entity target suppression, low-bit
stability, and generation/relearning audits. QR-
LoRA-style concentration addresses the low-bit sta-
bility component but can fail the semantic-neighbor
axis; AR-constrained LoRA restores that axis but

Table 24: Proxy-only post-training extension on Llama-3.2-3B (seed 42). The jog data are limited to instruction- and preference-style proxies at 5/10/20/50 steps, so this table should be read as a focused extension of Table 22, not as a replacement for the broader multi-jog audit.

Method	After TGT	Worst Proxy Jog	Worst mode	Max resurfaced	Worst N/G
Q-ROU	13/13	12/13	instr. proxy / 50	1	8/8 / 5/7
Q-ROU+PTI	13/13	13/13	instr. proxy / 50	0	8/8 / 5/7

Table 25: Broad post-training recurrence audit on Llama-3.2-3B (three seeds). This extension expands the jog family to neighbor, general, mixed, instruction proxy, and preference proxy at 5/10/20/50 steps. The PTI gain remains real, though smaller than in the narrower recurrence audit.

Method	Post-TGT	Worst broad	Worst case	Max mean	Worst proxy	Proxy worst
Q-ROU	11/13	6/13	mixed / 50	0.0339	9/13	instr. proxy / 50
Q-ROU+PTI	13/13	8/13	mixed / 50	0.0096	13/13	pref. proxy / 50

2209 under-forgets; memory-jogging shows that even
 2210 strong immediate deletion should be accompanied
 2211 by recurrence reporting; and PTI hardening demon-
 2212 strates that bounded recurrence robustness can be
 2213 improved without abandoning the original Q-ROU
 2214 deletion-and-retention operating point.

2215 **Step Sweep on Expanded 65-Probe Set.** Ta-
 2216 ble 27 extends the step-Pareto analysis to the ex-
 2217 panded evaluation set. At 20 steps, both selective
 2218 methods are still under-saturated, with GA+AR
 2219 slightly ahead on target removal (11/25 vs. 8/25).
 2220 By 40 steps, Q-ROU reaches full 25/25 target sup-
 2221 pression while GA+AR remains at 13/25. At 60
 2222 steps, GA+AR partially closes the target gap (22/25)
 2223 but still trails Q-ROU, whose neighbor retention
 2224 remains stable (17/20 FP16) and whose general be-
 2225 havior stays intact throughout.

2226 **TOFU Multi-Split Comparison.** Table 28 evalu-
 2227 ates both methods across all three TOFU subsets. Q-
 2228 ROU outperforms GA+AR on forgetting and leak-
 2229 age across all three subsets under the fixed 160-
 2230 step budget, but the hardest forget10 row should
 2231 be read as a transitional fixed-budget checkpoint
 2232 rather than as a stable operating ceiling. The for-
 2233 get05/forget10 values come from final hard-split
 2234 evaluations with explicit retain splits; long-horizon
 2235 sweeps are used solely for interpretation. A three-
 2236 seed confirmation on the tuned hard-split operating
 2237 points (240 steps, lf40_s240) gives mean results
 2238 of 99.83% Forget@0.01, 1.67% leakage, 98.33%
 2239 utility, and 94.67% truth on forget05, and 94.42%
 2240 Forget@0.01, 4.92% leakage, 90.83% utility, and
 2241 87.33% truth on forget10. We do not fold these
 2242 tuned rows into the fixed-budget main table, but
 2243 they show that the tuned hard-split trade-off is not

a single-seed artifact.

2244
 2245 As a smaller-model calibration check, we also
 2246 evaluated the TOFU matrix on Llama-3.2-3B with
 2247 explicit retain splits. Table 29 reports the resulting
 2248 rows and additional operating-point sweeps. We
 2249 treat the 80-step rows as *matched-budget diagnos-*
 2250 *tics*: they show that Q-ROU is decisive on forget01,
 2251 clearly ahead on forget05 forgetting/leakage but not
 2252 yet saturated, and not a successful operating point
 2253 on forget10. The additional Q-ROU sweep rows
 2254 then make the tuning effort explicit by sweeping
 2255 only Q-ROU while holding the matched GA+AR
 2256 reference fixed, rather than silently replacing the
 2257 harder rows with a hand-picked checkpoint. On
 2258 forget05, the sweep converts the under-saturated
 2259 80-step row into near-complete or complete dele-
 2260 tion with high utility. On forget10, additional effort
 2261 sharply reduces forget probability and leakage, but
 2262 the cleanest hard-split point remains utility-limited
 2263 (94.25% Forget@0.01 and 89.5% utility at 240
 2264 steps), while the 400-step row slightly improves
 2265 pass count but worsens tail probability.

2266 Because thresholded pass counts can be misread
 2267 when shown alone, we also report the same raw
 2268 forget probabilities under multiple thresholds in Ta-
 2269 ble 30. This table is meant to prevent the reader
 2270 from over-interpreting a single $P < 0.01$ cut: for ex-
 2271 ample, the under-saturated forget10 80-step row is
 2272 weak at strict thresholds but already shows a lower
 2273 probability mass than GA+AR, while the additional
 2274 Q-ROU sweep rows show the operating-point path
 2275 from mass reduction to strict deletion. They also ex-
 2276 pose a real small-model utility boundary: stronger
 2277 Q-ROU settings can recover most hard-split forget
 2278 cases, but not without lowering retain utility relative
 2279 to the matched GA+AR reference.

Table 26: Deployed-NF4 evaluation on the harshest 50-step slices from the broad Llama PTI audit. We reevaluate the ‘neighbor’, ‘general’, ‘mixed’, and ‘instruction proxy’ checkpoints under real BnB NF4. The PTI ordering is preserved and slightly widened on this subset.

Method	Source TGT	NF4 broad	NF4 worst case	NF4 proxy	Proxy worst	NF4 path
Q-ROU	11/13	5/13	mixed / 50	9/13	instr. proxy / 50	BnB NF4
Q-ROU+PTI	13/13	10/13	mixed / 50	13/13	instr. proxy / 50	BnB NF4

Table 27: Step sweep on the 65-probe expanded set (Llama-3.2-3B). Q-ROU is under-suppressed at 20 steps, cleanly overtakes GA+AR by 40 steps, and remains ahead at 60 steps while maintaining stable neighbor retention.

Steps	Q-ROU				GA+AR			
	FP TGT	FP NBR	I4 TGT	Δ_z	FP TGT	FP NBR	I4 TGT	Δ_z
20	8/25	20/20	10/25	-0.022	11/25	20/20	12/25	-0.030
40	25/25	17/20	25/25	-0.015	13/25	19/20	13/25	-0.028
60	25/25	17/20	25/25	-0.014	22/25	17/20	21/25	-0.016

Table 28: TOFU multi-split evaluation (Llama-3.1-8B, fixed 160-step budget, final hard-split evaluations for forget05/forget10).

Split	Method	Forget \bar{P}	Gen Leak	Forget@0.01	n
forget01	Q-ROU	0.00003	0.0%	100.0%	40
	GA+AR	0.22023	10.0%	52.5%	40
forget05	Q-ROU	0.00053	0.0%	100.0%	200
	GA+AR	0.41320	34.5%	11.0%	200
forget10	Q-ROU	0.05109	7.0%	72.0%	400
	GA+AR	0.45833	36.75%	5.75%	400

2280 Adding the other baselines from the same final
2281 hard-split evaluation clarifies the trade-off structure.
2282 At fixed 160 steps on 8B forget10, RMU achieves
2283 stronger raw forgetting than Q-ROU (Forget@0.01
2284 = 85.25% vs. 72.0%; leakage = 1.25% vs. 7.0%),
2285 but at the cost of severe utility/truth collapse (Util-
2286 ity = 12.0%, Truth = 66.67%). This is why the
2287 fixed-budget forget10 row is not, by itself, a clean
2288 all-baseline dominance result. The long-horizon
2289 sweeps below resolve that ambiguity.

2290 C.4 Forget10 Saturation Trajectory

2291 The main-text TOFU table reports the final 8B fixed-
2292 160 checkpoint for forget10 (72.0% Forget@0.01
2293 in the final hard-split evaluation), but this fixed-
2294 budget number is not stable enough to be interpreted
2295 as the method’s ceiling. Across standalone fixed-
2296 160 Full-Q-ROU-style runs on the same 8B split,
2297 Forget@0.01 spans 48.25% \rightarrow 77.5% while utility
2298 remains relatively high. We therefore extended the
2299 hardest 8B case directly rather than relying solely

on the 3B precedent. 2300

Figure 4 plots the 8B long-horizon step sweeps 2301
on Llama-3.1-8B: two Q-ROU seeds plus matched 2302
RMU and GA+AR comparator runs, all evaluated 2303
every 32 steps to step 640. These results indicate 2304
the performance trajectory clearly. Q-ROU 2305
is materially under-saturated at step 160 on both 2306
seeds (seed 42: 40.75%; seed 43: 47.50% For- 2307
get@0.01), then enters a high-performing regime 2308
between roughly steps 224 and 416. Seed 42 2309
first reaches near-complete forgetting at step 288 2310
(99.25%, Utility 93.0%, Leak 1.25%) and attains 2311
its best recorded score at step 384 (99.50%, Utility 2312
97.0%, Leak 0.5%). Seed 43 reaches 100% For- 2313
get@0.01 by step 352 and attains its best recorded 2314
score at step 512 (100%, Utility 92.0%, Leak 2315
0.0%). 2316

The matched long-horizon comparators 2317
strengthen this interpretation. RMU peaks early 2318
(best score at step 96 with Forget@0.01 = 91.0%) 2319
but preserves only 5.0% utility and never exceeds 2320
16.0% utility beyond step 128, degrading to 2321
Forget@0.01 = 70.0% and Leak = 5.25% by 2322
step 640. GA+AR never escapes its low-forgetting 2323
regime: across two independent matched-horizon 2324
runs, it peaks at only 6–10% Forget@0.01 around 2325
step 160 with leakage above 33% and remains there 2326
through step 640. The main lesson is therefore 2327
not “forget10 is hard at 160 steps”, but rather 2328
“the standard 160-step budget is under-saturated 2329
on 8B, and Q-ROU reaches a practically useful 2330
saturated regime well before the terminal 640-step 2331
checkpoint.” 2332

In short, the old 3B rescue argument is no longer 2333
the primary evidence. The key evidence is now 2334
direct 8B long-horizon confirmation: the standard 2335
160-step protocol underestimates Q-ROU on for- 2336
get10, the high-performing regime appears substan- 2337
tially before the terminal 640-step checkpoint, and 2338
the exact best checkpoint is seed-sensitive within 2339
that saturated region. 2340

**TOFU-Small Cross-Model Transfer Screen (6 2341
Methods, 2 Models).** We evaluated the TOFU 2342

Subset	Operating point	Method	Steps	λ_f	Forget@0.01	Gen leak	Utility	Truth
forget01	GA+AR ref	ga_ar	80	–	0.0%	35.0%	100.0%	96.0%
forget01	s80_lf20	qrou	80	20.0	100.0%	2.5%	99.5%	92.0%
forget05	GA+AR ref	ga_ar	80	–	3.5%	34.5%	100.0%	90.0%
forget05	s80_lf20	qrou	80	20.0	54.5%	11.0%	97.5%	90.0%
forget05	lf40_s160	qrou	160	40.0	98.0%	3.5%	99.0%	92.0%
forget05	lf80_s160_ar60	qrou	160	80.0	99.0%	4.0%	98.0%	90.0%
forget05	lf40_s240	qrou	240	40.0	100.0%	2.0%	99.0%	92.0%
forget05	lf40_s400	qrou	400	40.0	100.0%	2.5%	99.0%	86.0%
forget10	GA+AR ref	ga_ar	80	–	4.0%	34.0%	100.0%	98.0%
forget10	s80_lf20	qrou	80	20.0	2.2%	25.2%	99.0%	96.0%
forget10	lf40_s160	qrou	160	40.0	67.2%	8.0%	87.5%	94.0%
forget10	lf80_s160_ar60	qrou	160	80.0	77.8%	9.0%	80.0%	90.0%
forget10	lf40_s240	qrou	240	40.0	94.2%	3.8%	89.5%	92.0%
forget10	lf40_s400	qrou	400	40.0	95.5%	3.0%	89.5%	94.0%

Table 29: Extended TOFU operating-point matrix on Llama-3.2-3B. The additional Q-ROU sweep rows reuse the matched GA+AR reference and probe harder forget subsets without silently replacing the baseline row.

Subset	Operating point	Method	$P < .005$	$P < .01$	$P < .02$	$P < .05$	Mean P	Max P
forget01	GA+AR ref	ga_ar	0.0%	0.0%	0.0%	0.0%	0.4174	0.9967
forget01	s80_lf20	qrou	97.5%	100.0%	100.0%	100.0%	0.0004	0.0061
forget05	GA+AR ref	ga_ar	3.5%	3.5%	5.0%	8.5%	0.3827	0.9846
forget05	s80_lf20	qrou	39.5%	54.5%	66.0%	72.5%	0.0585	0.5334
forget05	lf40_s160	qrou	97.5%	98.0%	98.0%	99.0%	0.0014	0.1115
forget05	lf80_s160_ar60	qrou	99.0%	99.0%	100.0%	100.0%	0.0003	0.0199
forget05	lf40_s240	qrou	100.0%	100.0%	100.0%	100.0%	0.0001	0.0020
forget05	lf40_s400	qrou	100.0%	100.0%	100.0%	100.0%	0.0001	0.0003
forget10	GA+AR ref	ga_ar	2.5%	4.0%	4.8%	5.0%	0.3999	0.9912
forget10	s80_lf20	qrou	1.2%	2.2%	4.5%	18.2%	0.1987	0.8167
forget10	lf40_s160	qrou	64.2%	67.2%	72.0%	77.0%	0.0453	0.4640
forget10	lf80_s160_ar60	qrou	76.2%	77.8%	78.2%	79.2%	0.0624	0.7832
forget10	lf40_s240	qrou	92.8%	94.2%	94.5%	95.0%	0.0108	0.3082
forget10	lf40_s400	qrou	95.2%	95.5%	95.5%	95.8%	0.0142	0.5552

Table 30: TOFU threshold sensitivity on Llama-3.2-3B. The table reports the same raw forget probabilities at multiple pass thresholds, making the operating-point dependence visible rather than relying only on $P < 0.01$.

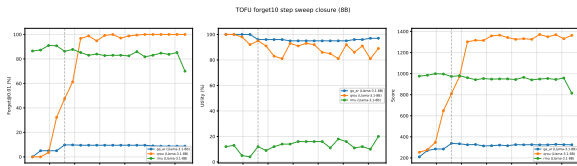


Figure 4: 8B forget10 long-horizon step sweeps on Llama-3.1-8B (640 steps, TOFU aggressive profile). Two Q-ROU seeds show that the standard 160-step budget is under-saturated and that a high-performing saturated regime appears by roughly steps 288–416, while RMU peaks early with severe utility collapse and GA+AR never exits its low-forgetting regime. The vertical dashed line indicates the default 160-step budget used in Table 1.

2343 transfer screen on the final robust forget01 setup
 2344 over two architectures (Llama-3.2-3B and Llama-
 2345 3.1-8B) and six methods (Q-ROU, GA, NPO,
 2346 GA+AR, RMU, RepBend), using a matched high-

budget setting (160 steps for both Q-ROU and the
 comparison methods) with step-trace logging ev-
 ery 32 steps. To avoid cherry-picking, Table 32
 reports the fixed terminal checkpoint (step 160) for
 all methods; step-trace best checkpoints are pro-
 vided separately as diagnostic evidence.

Evaluation Metrics. The robust proxy evaluation
 utilizes the following definitions:

- **Forget@ P :** The percentage of forget queries
 where the model’s mean per-token probability
 for the correct answer falls below threshold
 P (e.g., 0.01 or 0.05). For context, the base-
 line Llama-3.2-3B and Llama-3.1-8B models
 assign mean probabilities of roughly 16.5%
 and 17.9% to the target answers prior to un-
 learning, making $P = 0.05$ a clear degradation
 and $P = 0.01$ a strict drop of over an order
 of magnitude representing near-complete sup-

Table 31: Selected 8B forget10 long-horizon checkpoints (Llama-3.1-8B, aggressive profile).

Method	Seed	Step	F@0.01	F@0.05	Leak	Utility	Truth
Q-ROU	42	160	40.75%	51.75%	11.50%	96.0%	100.0%
Q-ROU	42	288	99.25%	100.0%	1.25%	93.0%	95.83%
Q-ROU	42	384	99.50%	99.75%	0.50%	97.0%	95.83%
Q-ROU	42	640	100.0%	100.0%	0.00%	82.0%	95.83%
Q-ROU	43	160	47.50%	69.50%	11.25%	95.0%	95.83%
Q-ROU	43	352	100.0%	100.0%	1.00%	93.0%	87.50%
Q-ROU	43	512	100.0%	100.0%	0.00%	92.0%	91.67%
Q-ROU	43	640	100.0%	100.0%	0.00%	89.0%	91.67%
RMU	42	96	91.00%	94.75%	1.25%	5.0%	66.67%
RMU	42	160	86.25%	90.75%	0.75%	12.0%	66.67%
RMU	42	640	70.00%	79.50%	5.25%	20.0%	66.67%
GA+AR	42	160	9.75%	9.75%	33.75%	96.0%	91.67%
GA+AR	42	640	8.75%	9.00%	35.00%	97.0%	91.67%
<i>GA+AR Run 2 (independent rerun)</i>							
GA+AR	42	160	6.00%	6.00%	39.75%	100.0%	100.0%
GA+AR	42	640	5.75%	6.00%	36.75%	99.0%	91.67%

2365 pression on the tested probes.

2366 • **Leak:** The percentage of forget queries where
2367 unconstrained greedy text generation produces
2368 the exact answer string or > 50% of its mean-
2369 ingful words.

2370 • **Utility:** The percentage of retained queries
2371 where the answer probability remains at or
2372 above 30% of the unmodified baseline model’s
2373 probability.

2374 • **Truth:** The percentage of general knowledge
2375 queries where the sequence log-probability
2376 of the true answer strictly exceeds that of a
2377 randomly mismatched false answer.

2378 All 12 model×method runs completed without
2379 runtime failures in this sweep. Under fixed step-
2380 160 comparison, Q-ROU remains top-ranked on
2381 both models (1406.51 on 3B, 1401.52 on 8B), with
2382 strict/relaxed forgetting both at 100%, zero leakage,
2383 and high utility (98–100%). GA reaches 100% for-
2384 getting with zero leakage but collapses utility to 0%
2385 on both models; RMU maintains high forgetting but
2386 with low utility (8% on 3B, 30% on 8B); GA+AR
2387 preserves high utility but reaches only 50% forget-
2388 ting with non-zero leakage (5.0% on 3B, 7.5% on
2389 8B) and much higher compute cost.

2390 For score interpretation, the benchmark uses

$$\begin{aligned} \text{Score} = & 1000 \cdot \text{Forget} - 450 \cdot \text{Leak} \\ & + 260 \cdot \text{Utility} + 120 \cdot \text{Truth} \\ & + 40 \cdot \min(\text{RetainRatio}, 1.5), \end{aligned}$$

2391

where RetainRatio is the mean retain probability
ratio used in the evaluation pipeline; this weight-
ing heavily rewards forgetting. For readability,
Table 32 reports Forget/Leak/Utility/Truth in %,
but the score is computed from their unit-interval
forms (0–1), consistent with the scoring rule used
throughout the TOFU evaluation. Using practical
guardrails (utility $\geq 80\%$, leak $\leq 5\%$), Q-ROU is
the only method that remains feasible on both mod-
els in this run. Because this is still a single-seed
screen (seed 42), we treat it as structured diagnostic
evidence rather than a final robustness estimate.

Table 33 shows heterogeneous best-step locations
across methods and models. Q-ROU peaks at step
128 on both models, whereas GA and NPO peak
later on 8B, and RMU/RepBend peak earlier (step
64) on both models. This pattern further supports
reporting fixed-budget terminal checkpoints for pri-
mary comparisons and treating per-method best-
step traces as secondary diagnostics.

C.5 Mechanism Isolation: Active Retention Transplant

To isolate the contribution of Active Retention,
we conduct a transplant experiment: adding AR
($\lambda_a = 80$, KL-divergence on pre-computed neigh-
bor logits) to each baseline, and removing it from
Q-ROU. Table 34 presents the full results including
both FP and INT4 evaluations.

The results establish three findings. First, **AR is necessary for selective forgetting in Q-ROU:** removing AR causes neighbor retention to collapse from 8/8 to 2/8 (FP) and 0/8 (INT4), while target removal accelerates from 4/13 to 13/13 at 20 steps. This confirms that AR acts as a precision brake—without it, Q-ROU’s forget loss overpowers the remaining regularization components (EWC, orthogonality, margin), producing aggressive but non-selective forgetting qualitatively similar to the baselines. The 2/8 residual neighbor retention (vs 0/8 for baselines without AR) reflects the partial protection from SLUG and EWC, but this is insufficient for reliable deployment.

Second, **AR is sufficient to prevent neighbor collapse in baselines:** adding it to all three baselines restores neighbor retention from 0/8 to 8/8 and general knowledge from 0–5/7 to 7/7, even without EWC, SLUG, orthogonality, or QuantNoise. However, this sufficiency is *neighbor-metric scoped*: it does not by itself resolve the token-generation tradeoff or semantic leakage patterns that emerge under extended training and adversarial evaluation

Table 32: TOFU-small cross-model transfer summary (2 models, 6 methods, robust profile, matched 160/160 budget). Values are reported at fixed final step 160; mean \pm sd across Llama-3.2-3B and Llama-3.1-8B.

Method	Forget@0.01 (%)	Forget@0.05 (%)	Leak (%)	Utility (%)	Truth (%)	Score	Train+Eval (s)
Q-ROU	100.00\pm0.00	100.00\pm0.00	0.00\pm0.00	99.00\pm1.00	97.91\pm2.08	1404.02\pm2.49	231.20 \pm 45.59
GA	100.00 \pm 0.00	100.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	27.08 \pm 6.25	1032.50 \pm 7.50	213.78\pm54.40
NPO	35.00 \pm 17.50	62.50 \pm 30.00	0.00 \pm 0.00	46.00 \pm 35.00	93.75 \pm 6.25	593.77 \pm 70.13	232.62 \pm 50.92
GA+AR	50.00 \pm 0.00	50.00 \pm 0.00	6.25 \pm 1.25	98.00 \pm 1.00	97.91 \pm 2.08	880.20 \pm 11.89	678.04 \pm 103.25
RMU	96.25 \pm 3.75	98.75 \pm 1.25	0.00 \pm 0.00	19.00 \pm 11.00	68.75 \pm 6.25	1099.65 \pm 1.17	228.09 \pm 59.95
RepBend	60.00 \pm 0.00	63.75 \pm 3.75	0.00 \pm 0.00	68.00 \pm 2.00	77.09 \pm 2.08	889.55 \pm 3.39	225.55 \pm 63.88

Table 33: Step-trace best checkpoints by method (32-step interval; checkpoints at 32/64/96/128/160). Final cross-method comparison in Table 32 remains fixed at step 160 for fairness.

Method	3B Best Step	3B Best Score	3B Score@160	8B Best Step	8B Best Score	8B Score@160
Q-ROU	128	1411.46	1406.51	128	1401.88	1401.52
GA	32	1030.00	1025.00	160	1040.00	1040.00
NPO	128	526.24	523.64	160	663.90	663.90
GA+AR	160	892.08	892.08	128	868.55	868.31
RMU	64	1098.76	1098.48	64	1106.57	1100.82
RepBend	64	894.73	886.16	64	898.30	892.94

Table 34: Mechanism isolation via Active Retention transplant (3B multi-entity, 20 steps, single-GPU execution). Removing AR from Q-ROU causes NBR collapse to 2/8 (FP) and 0/8 (INT4); adding AR to baselines recovers NBR from 0/8 to 8/8. INT4 target forgetting often matches or exceeds FP while retention remains locked.

Method	AR	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z
GA	\times	13/13	0/8	0/7	13/13	0/8	0/7	+0.000
GA+AR	\checkmark	4/13	8/8	7/7	6/13	8/8	7/7	-0.021
GD	\times	13/13	0/8	4/7	13/13	0/8	4/7	+0.000
GD+AR	\checkmark	5/13	8/8	7/7	7/13	8/8	7/7	-0.018
RB	\times	13/13	0/8	5/7	13/13	0/8	5/7	+0.000
RB+AR	\checkmark	5/13	8/8	7/7	7/13	8/8	7/7	-0.015
Q-ROU	\checkmark	4/13	8/8	7/7	4/13	8/8	7/7	-0.021
Q-ROU -AR	\times	13/13	2/8	7/7	13/13	0/8	7/7	+0.000

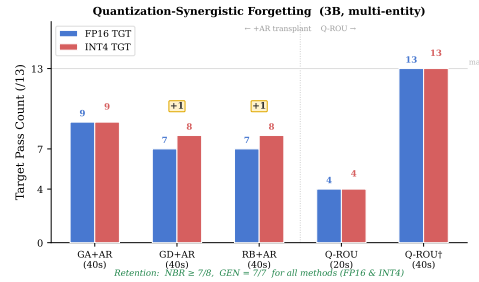


Figure 5: Low-bit target drift on the 3B multi-entity audit. AR-equipped baselines show small target-side INT4 gains while maintaining retention, whereas converged Q-ROU reaches the same complete target suppression in FP16 and INT4.

2443 (Section D.4; Section D.3).

2444 Third, and perhaps most importantly, **INT4**
 2445 **quantization can selectively enhance target forgetting**
 2446 **when AR is present**. Fig. 5 visualizes this
 2447 pattern across all +AR baselines (40 steps) along-
 2448 side Q-ROU at both 20 and 40 steps: comparing FP
 2449 TGT and I4 TGT bars reveals consistent enhance-
 2450 ment (+0 to +1 cases) under AR while retention
 2451 remains stable, and Q-ROU[†] (40 steps) achieving
 2452 complete target removal (13/13) in both precisions.

2453 This *quantization-synergistic forgetting* arises be-
 2454 cause AR constrains the parameter configuration
 2455 to a structurally stable region for retention. The
 2456 small perturbations from INT4 rounding can then
 2457 only meaningfully affect already-borderline target
 2458 probabilities, nudging them below the $P < 0.01$
 2459 removal threshold without disturbing the firmly-

2460 anchored neighbor outputs. Full Q-ROU at 20 steps
 2461 shows matched FP/INT4 target counts (4/13) be-
 2462 cause QuantNoise has already flattened the loss
 2463 landscape with respect to quantization perturba-
 2464 tions, leaving no borderline cases for INT4 to push.
 2465 At 40 steps, Q-ROU achieves 13/13 target in both
 2466 precisions (fully converged, no borderline cases)
 2467 with near-zero zombie delta ($\Delta_z \approx 0$), confirming
 2468 that QuantNoise creates near-identical FP/INT4 be-
 2469 havior throughout the optimization trajectory.

2470 The +AR methods show reduced target removal
 2471 at 20 steps (4–5/13 vs 13/13 without AR), reflecting
 2472 AR’s deliberate brake on forgetting speed.

2473 To verify convergence, we extend the transplant
 2474 to 40 steps (Table 18). At this duration, GA+AR
 2475 and GD+AR each reach 7/13 TGT in both FP and
 2476 INT4 (with 8/8 NBR and 7/7 GEN), while RB+AR

2477 reaches 7/13 in FP and 8/13 in INT4 with the same
 2478 8/8 NBR and 7/7 GEN pattern. Q-ROU at 40 steps
 2479 surpasses all +AR baselines by 5 target cases against
 2480 the strongest baseline (13/13 vs 8/13) while trad-
 2481 ing at most one neighbor case under INT4. This
 2482 demonstrates that Q-ROU’s additional components
 2483 (KL-to-uniform forget loss, QuantNoise, orthogo-
 2484 nality, margin) provide a substantial advantage be-
 2485 yond what AR alone contributes. The quantization-
 2486 synergistic pattern persists in +AR baselines at +0 to
 2487 +1 case under INT4 while retention remains locked.
 2488 As a control, GA *without* AR at 40 steps destroys not
 2489 only all neighbor knowledge (0/8) but also all gen-
 2490 eral knowledge (0/7)—every evaluation prompt pro-
 2491 duces near-zero probability, confirming total model
 2492 collapse. This 40-step control eliminates the alter-
 2493 native explanation that baselines fail merely due to
 2494 insufficient training: extended training *worsens* the
 2495 collapse when AR is absent, while AR-equipped
 2496 methods converge to strong forgetting with intact
 2497 retention.

2498 C.6 Ablation Study

2499 The main text already summarizes the 3B compo-
 2500 nent ablation covering the core components (AR,
 2501 QuantNoise, Orthogonality, EWC, Margin). Here,
 2502 we provide the 8B scale ablation on TOFU and the
 2503 adaptive mechanism ablation.

2504 **8B Scale Ablation on TOFU.** At 8B (Table 35:
 2505 TOFU forget01, 160 steps), all six ablation vari-
 2506 ants achieve identical Forget@0.01 = 100.0% with
 2507 $|\Delta S| \leq 2.54$. Margin and EWC removal produce
 2508 exactly zero measurable effect on superficial pass
 2509 counts, confirming that parameter redundancy han-
 2510 dles capacity demands at 8B. However, these auxil-
 2511 iary losses act as a geometric safety net (defense-in-
 2512 depth): for instance, the Orthogonality constraint
 2513 implicitly bounds the rotation of general-knowledge
 2514 subspaces (Appendix Theorem 2). Thus, while
 2515 they may appear redundant on easy fixed-budget
 2516 checkpoints, they provide geometric stability that
 2517 becomes relevant under harder regimes and longer
 2518 horizons (Appendix C.4; Appendix C.16). How-
 2519 ever, on harder forget10 at longer horizons, utility
 2520 ordering becomes seed-sensitive (Appendix C.6).

2521 **Adaptive Mechanism Ablation.** Table 36 reports
 2522 the four-way ablation of adaptive mechanisms on
 2523 the 65-probe expanded set.

2524 We emphasize that this ablation is intentionally
 2525 reported at fixed-budget binary endpoints (T/N/G

Table 35: 8B component ablation on TOFU forget01 (Llama-3.1-8B, 160 steps, seed=42). All variants achieve identical Forget@0.01 = 100%; score variation is ≤ 2.54 .

Configuration	F@0.01	F@0.05	Leak	Utility	Truth	Score
Full Q-ROU	100.0%	100.0%	0.0%	100.0%	91.7%	1396.01
–Orthogonality	100.0%	100.0%	0.0%	99.0%	91.7%	1393.47
–Margin	100.0%	100.0%	0.0%	100.0%	91.7%	1396.01
–EWC	100.0%	100.0%	0.0%	100.0%	91.7%	1396.20
–Ortho–Margin	100.0%	100.0%	0.0%	99.0%	91.7%	1393.47
–Ortho–Margin–EWC	100.0%	100.0%	0.0%	100.0%	91.7%	1395.94

pass counts). Adaptive variants are primarily de- 2526
 signed to improve robustness properties that are 2527
 only weakly reflected in these endpoints (e.g., 2528
 threshold sensitivity and precision brittleness), and 2529
 should therefore be interpreted alongside the thresh- 2530
 old robustness analyses in Section C.12. 2531

For the adaptive-mechanism subtable (Standard 2532
 / Adaptive-QN / Adaptive-AR / Full Adaptive), we 2533
 verified the exported LaTeX rows against the under- 2534
 lying count totals (25/20/20 in FP and INT4, plus 2535
 legacy 13/8/7 regression metadata for qrou_40s). 2536

3B Component-Level Ablation Details. Ta- 2537
 ble 37 supplements the main-text discussion of 2538
 scale-dependent component contribution with con- 2539
 tinuous mean-probability values on Llama-3.2-3B 2540
 (20 steps, single-entity evaluation). All four config- 2541
 urations produce *identical* pass counts (FP16: TGT 2542
 2/7, NBR 9/9, GEN 7/7; INT4: TGT 3/7, NBR 2543
 9/9, GEN 7/7), confirming the same parameter- 2544
 redundancy reading summarized in the main text. 2545
 At the \bar{p}_{TGT} level, removing EWC actually *de-* 2546
creases target probability (0.012 vs. 0.025 for Full), 2547
 because EWC constrains parameter drift away from 2548
 the pre-training manifold, partially counteracting 2549
 the forget objective. The –Margin configuration 2550
 is *numerically identical* to Full Q-ROU (matching 2551
 to six decimal places), indicating that at 3B scale, 2552
 the margin-based separation loss contributes zero 2553
 incremental effect within the SLUG-constrained 2554
 parameter space. 2555

8B Component-Level Ablation on TOFU. Ta- 2556
 ble 38 extends the component ablation to Llama- 2557
 3.1-8B using the standardized TOFU forget01 2558
 benchmark (40 examples, 160 steps, seed=42). 2559
 Six configurations are tested: Full Q-ROU, and 2560
 five progressive removals of auxiliary compo- 2561
 nents (–Ortho, –Margin, –EWC, –Ortho–Margin, 2562
 –Ortho–Margin–EWC). 2563

All configurations achieve *identical* Forget@0.01 2564

Table 36: Adaptive mechanism ablation on the 65-probe expanded set (Llama-3.2-3B, 40 steps). Pass-count differences are small, positioning adaptive AR/QN primarily as stability instrumentation.

Configuration	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z
Standard (Fixed AR, Fixed QN)	25/25	17/20	20/20	25/25	16/20	18/20	-0.018
Adaptive QN Only ($\beta = 0.4$)	25/25	18/20	20/20	25/25	17/20	18/20	-0.017
Adaptive AR Only (max $\lambda_a = 80$)	25/25	17/20	20/20	25/25	16/20	18/20	-0.018
Full Adaptive (diagnostic variant)	24/25	17/20	20/20	25/25	16/20	18/20	-0.016

Table 37: 3B component ablation: continuous mean-probability values (Llama-3.2-3B, 20 steps, single-entity). All configurations yield identical binary pass counts.

Configuration	FP \bar{p}_{TGT}	FP \bar{p}_{NBR}	I4 \bar{p}_{TGT}	I4 \bar{p}_{NBR}	Δ_z
Full Q-ROU	0.025	0.937	0.023	0.921	-0.002
Q-ROU -Orthogonality	0.025	0.937	0.023	0.922	-0.002
Q-ROU -EWC	0.012	0.922	0.011	0.901	-0.001
Q-ROU -Margin	0.025	0.937	0.023	0.921	-0.002

2565 = 100.0%, Forget@0.05 = 100.0%, and zero generation leakage, confirming that at 8B scale, the
 2566 primary forget and retain objectives are sufficient
 2567 for complete target suppression on this split even
 2568 without auxiliary regularization.
 2569

2570 Step-curve analysis (5 checkpoints at 32-step intervals)
 2571 reveals that the Full configuration already
 2572 reaches 100% Forget@0.01 at step 96 on this split,
 2573 leaving no room for differentiation. The only measurable
 2574 difference is a 1% Utility drop ($\Delta = -1.0\%$)
 2575 when Orthogonality is removed, propagating to a
 2576 $\Delta S = -2.54$ score reduction. This is attributable to
 2577 Ortho’s ℓ_2 -regularization effect partially constraining
 2578 parameter drift away from the retain manifold.
 2579 Margin and EWC produce exactly zero measurable
 2580 effect at this scale ($\Delta S = 0.00$ and $+0.19$ respectively).
 2581

2582 The most aggressive variant
 2583 (-Ortho-Margin-EWC; effectively Q-ROU
 2584 reduced to forget + retain + AR + SLUG only)
 2585 achieves Score = 1395.94 versus Full’s 1396.01, a
 2586 Δ of just -0.08 . This provides additional evidence
 2587 for the parameter-redundancy hypothesis: at
 2588 8B scale with 40-example forget sets, the core
 2589 objectives alone carry all measurable performance.

2590 These values come from the same controlled 8B
 2591 easy-split ablation evaluations summarized in this
 2592 appendix.

2593 This easy-split redundancy story does not transfer
 2594 unchanged to harder TOFU settings. On 8B
 2595 forget10, the fixed-160 hard-split comparison initially
 2596 suggested a strong ordering Full > Core >

Core-no-QuantNoise, but the longer-horizon evaluations
 2597 refine that picture. Two independent 384-step
 2598 runs (Table 39) show that all three variants achieve
 2599 near-complete forgetting ($\geq 96.5\%$ Forget@0.01),
 2600 but the relative utility ordering is *not stable across*
 2601 *runs*: Run 1 places Full highest (Utility 95%) and
 2602 Core-no-QuantNoise lowest (87%), while Run 2
 2603 reverses this pattern (Core-no-QuantNoise 94%, Full
 2604 96%, Core 90%). The strongest surviving component
 2605 conclusion is therefore not about any single
 2606 auxiliary term, but rather that (i) all three variants
 2607 converge to near-complete forgetting at extended
 2608 horizon, and (ii) the relative utility ranking is seed-
 2609 sensitive, so single-run component claims should
 2610 be treated with caution on this harder split.
 2611

2612 C.7 Hyperparameter Sensitivity

2613 The Active Retention coefficient λ_a is the primary
 2614 hyperparameter controlling the trade-off between
 2615 forgetting aggressiveness and neighbor protection.
 2616 Table 40 reports Q-ROU performance on Llama-
 2617 3.2-3B multi-entity unlearning at 40 steps across
 2618 $\lambda_a \in \{40, 80, 120, 160\}$.

2619 The results reveal a smooth trade-off. Lower
 2620 λ_a values (40–80) permit aggressive forgetting
 2621 (13/13 in both FP and INT4) with one neighbor
 2622 case dropped (7/8) in both precisions. $\lambda_a = 120$
 2623 recovers that neighbor case *selectively* under INT4
 2624 (8/8 vs. its own FP score of 7/8), while maintaining
 2625 13/13 target removal. At $\lambda_a = 160$, both precisions
 2626 retain 8/8 neighbors while preserving 13/13 target
 2627 removal.

2628 Crucially, total pass counts remain at 27–28/28
 2629 across the entire four-fold range of λ_a , confirming
 2630 that the method does not require extensive hyper-
 2631 parameter tuning. We use $\lambda_a = 80$ as the default
 2632 throughout this paper for consistency across ex-
 2633 periments, while noting that $\lambda_a = 160$ maximizes
 2634 neighbor retention in this sweep without sacrificing
 2635 target removal.

Table 38: 8B component ablation on TOFU forget01 (Llama-3.1-8B, 160 steps, seed=42). All variants achieve identical perfect forgetting; score variation is ≤ 2.54 .

Configuration	F@0.01	F@0.05	Leak	Utility	Truth	Score	ΔS
Full Q-ROU	100.0%	100.0%	0.0%	100.0%	91.7%	1396.01	—
–Orthogonality	100.0%	100.0%	0.0%	99.0%	91.7%	1393.47	–2.54
–Margin	100.0%	100.0%	0.0%	100.0%	91.7%	1396.01	0.00
–EWC	100.0%	100.0%	0.0%	100.0%	91.7%	1396.20	+0.19
–Ortho–Margin	100.0%	100.0%	0.0%	99.0%	91.7%	1393.47	–2.54
–Ortho–Margin–EWC	100.0%	100.0%	0.0%	100.0%	91.7%	1395.94	–0.08

Table 39: 8B hard-split component evaluations on TOFU forget10 at a longer horizon (Llama-3.1-8B, 384 steps, seed=42). Two independent runs show that all variants achieve near-complete forgetting, but the relative utility ordering reverses across runs.

Run	Configuration	F@0.01	F@0.05	Leak	Utility	Truth	Score
1	Full Q-ROU	99.0%	100.0%	0.25%	95.0%	95.8%	1374.64
1	Core	100.0%	100.0%	0.25%	91.0%	95.8%	1373.63
1	Core-no-QuantNoise	100.0%	100.0%	0.0%	87.0%	95.8%	1362.45
2	Full Q-ROU	96.5%	99.5%	0.5%	96.0%	100.0%	1356.48
2	Core	99.5%	100.0%	0.25%	90.0%	91.7%	1361.70
2	Core-no-QuantNoise	99.25%	100.0%	0.25%	94.0%	91.7%	1369.53

Table 40: Sensitivity to λ_a (3B multi-entity, 40 steps). The method achieves $\geq 12/13$ target removal and $\geq 7/8$ neighbor retention across the full range, demonstrating robustness to hyperparameter choice.

λ_a	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z
40	13/13	7/8	7/7	13/13	7/8	7/7	≈ 0
80	13/13	7/8	7/7	13/13	7/8	7/7	≈ 0
120	13/13	7/8	7/7	13/13	8/8	7/7	≈ 0
160	13/13	8/8	7/7	13/13	8/8	7/7	≈ 0

2636 C.8 Computational Cost

2637 The computational cost and efficiency metrics of
 2638 Q-ROU relative to baselines are contextualized in
 2639 Section D. Table 41 reports measured wall-clock
 2640 costs for the larger same-work probability and gen-
 2641 eration evaluations. The fine-grained AR oper-
 2642 ating point roughly doubles the retain-anchor set
 2643 and therefore costs more than default Q-ROU, but
 2644 the observed per-seed training times remain in the
 2645 one-to-two-minute range on the logged single-GPU
 2646 runs used for that table. Table 42 then places
 2647 Q-ROU, `qrou_fg_ar`, GA+AR, and `gaar_fg_ar`
 2648 on the same 81-probe same-work benchmark and
 2649 the same RTX 3090 hardware slice. At 80 steps,
 2650 `qrou_fg_ar` is the most expensive operating point
 2651 (787 s train, 10.4 GB peak VRAM), but it is also the
 2652 only configuration in this comparison that preserves
 2653 full 13/13 target suppression while materially im-
 2654 proving same-work retention over default Q-ROU.

C.9 Expanded Evaluation

2655 Table 43 reports a focused five-method subset of
 2656 the expanded 65-probe evaluation on Llama-3.2-
 2657 3B (multi-entity, 40 steps), matching the central
 2658 selective-forgetting comparison discussed in the
 2659 main text. In this subset, Q-ROU reaches 25/25
 2660 target removal in FP16 and remains at 25/25 under
 2661 INT4 while preserving substantially more neigh-
 2662 bor/general behavior than aggressive baselines. In
 2663 contrast, GA and RepBend achieve 100% target
 2664 removal only at the cost of complete neighbor col-
 2665 lapse (0/20). Because this 40-step row is used
 2666 as a central selective-forgetting result, Table 44
 2667 adds a focused three-seed confirmation for the two
 2668 strongest selective methods at the same operating
 2669 point.
 2670

C.10 Depth Probing and Adversarial Robustness

2671 Standard evaluation probes primarily test direct
 2672 knowledge recall. To assess whether Q-ROU
 2673 achieves deeper, conceptual-level knowledge sup-
 2674 pression, we design six structured probing protocols
 2675 with 55 probes total.
 2676

2677 *Protocol A* (reverse association, 12 probes) tests
 2678 **one-step reverse lookup**: given a distinguishing at-
 2679 tribute or description of the forgotten entity, can the
 2680 model retrieve the entity name itself? For example,
 2681 “The fictional school with moving staircases is” \rightarrow
 2682 “Hogwarts” inverts the standard entity \rightarrow attribute
 2683

Table 41: Runtime accounting for the larger same-work audits. Values are mean \pm standard deviation over seeds.

Track	Method	Seeds	Train s	Wall s	Retain texts
custom neighbor audit	GA+AR	6	26.8 \pm 0.5	26.9 \pm 0.5	25.0
custom neighbor audit	GA+FG-AR	6	43.0 \pm 0.4	43.2 \pm 0.4	57.0
custom neighbor audit	Q-ROU	6	47.2 \pm 0.2	47.9 \pm 0.2	25.0
custom neighbor audit	Q-ROU+FG-AR	6	85.7 \pm 0.4	86.7 \pm 0.4	57.0
samework generation audit	GA+AR	3	37.7 \pm 0.6	37.9 \pm 0.7	25.0
samework generation audit	GA+FG-AR	3	60.1 \pm 0.8	60.3 \pm 0.8	57.0
samework generation audit	Q-ROU	3	68.2 \pm 0.4	69.2 \pm 0.6	25.0
samework generation audit	Q-ROU+FG-AR	3	123.0 \pm 0.6	124.2 \pm 0.6	57.0

Table 42: Uniform runtime accounting on the larger same-work benchmark (81 probes, single RTX 3090, one seed). This table uses the same hardware and layer set for all four methods, making relative training cost directly comparable.

Steps	Method	Train s	Wall s	Peak VRAM MB	TGT	SW	ID	GEN
40	Q-ROU	213.3	222.5	9853.8	13/13	4/24	12/16	13/16
40	Q-ROU+FG-AR	384.3	390.6	10389.5	12/13	14/24	15/16	16/16
40	GA+AR	91.1	94.7	8548.0	7/13	16/24	15/16	16/16
40	GA+FG-AR	175.4	181.7	8934.8	1/13	21/24	16/16	16/16
80	Q-ROU	434.2	442.9	9853.8	13/13	4/24	12/16	14/16
80	Q-ROU+FG-AR	787.1	797.2	10389.5	13/13	12/24	14/16	16/16
80	GA+AR	244.9	250.8	8548.0	12/13	5/24	11/16	15/16
80	GA+FG-AR	397.7	404.1	8934.8	5/13	22/24	16/16	16/16

Table 43: Focused five-method subset of the expanded 65-probe evaluation on Llama-3.2-3B (multi-entity, 40 steps), matching the central selective-forgetting comparison from Table 3. Q-ROU reaches 100% target removal in both FP16 and INT4 while maintaining substantially higher neighbor/general retention than aggressive baselines.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z
Baseline	5/25	20/20	20/20	5/25	18/20	20/20	-0.055
Q-ROU	25/25	17/20	20/20	25/25	16/20	18/20	-0.018
GA+AR	13/25	18/20	20/20	13/25	15/20	19/20	-0.026
GA	25/25	0/20	0/20	25/25	0/20	0/20	+0.000
RepBend	25/25	0/20	12/20	25/25	0/20	11/20	-0.013

2684 direction by querying attribute \rightarrow entity. This tests
 2685 whether the suppression of the entity’s name holds
 2686 when approached from the property side, rather
 2687 than by naming the entity directly.

2688 *Protocol B* (multistep reasoning, 10 probes) tests
 2689 **multi-hop inference chains**: the probe does not
 2690 name the forgotten entity or its attributes directly
 2691 but instead requires the model to traverse two or
 2692 more inferential steps to reach the target fact. For ex-
 2693 ample, “The sport played on broomsticks in a mag-
 2694 ical school” \rightarrow “Quidditch” requires the model to
 2695 first identify the school (Harry Potter’s Hogwarts),
 2696 then retrieve the sport associated with it—neither
 2697 step alone is sufficient. This tests whether knowl-
 2698 edge that must be *inferred via intermediate concepts*

is also suppressed.

Protocol C (paraphrased extraction, 15 probes) rephrases the original prompt while preserving the same target answer, testing whether surface-level lexical changes bypass suppression. *Protocol D* (reconstruction, 9 probes) provides incremental hints to test resistance to guided reconstruction. *Protocol E* (in-context priming, 6 probes) provides related context explicitly in the prompt to prime knowledge retrieval. *Protocol F* (negation-based extraction, 3 probes) uses negated or contrastive framing to elicit the target answer.

Table 45 presents the results. Q-ROU achieves **100% suppression (55/55)** across all six protocols in both FP16 and INT4, with mean target probability below 0.002 in every protocol. This result substantially strengthens the evidence for conceptual-level suppression: the method suppresses not only direct recall but also reverse associations (Protocol A), multi-hop reasoning chains (Protocol B), hint-guided reconstruction (Protocol D), and all three adversarial extraction strategies.

GA+AR provides partial but incomplete defense (43/55 in both FP16 and INT4; 78%). Its critical weakness is *in-context priming*: when relevant context is explicitly provided in the prompt, GA+AR suppresses only 3/6 probes (50%) with a mean probability of 0.204—two orders of magni-

Table 44: Three-seed confirmation of the expanded 65-probe 40-step operating point (Llama-3.2-3B). Values are mean pass counts over seeds 42/43/44; parentheses show population standard deviation. This confirms that the main 40-step ordering is stable rather than a single-seed artifact.

Method	FP16				INT4			
	TGT/25	NBR/20	GEN/20	Δ_z	TGT/25	NBR/20	GEN/20	Δ_z
Q-ROU	25.0 (0.0)	17.0 (0.0)	19.7 (0.5)	-0.0181 (0.0006)	25.0 (0.0)	15.7 (0.5)	18.0 (0.0)	-0.0181 (0.0006)
GA+AR	13.0 (0.0)	19.0 (0.0)	20.0 (0.0)	-0.0281 (0.0000)	13.0 (0.0)	15.0 (0.0)	19.0 (0.0)	-0.0281 (0.0000)

Table 45: Depth probing on Llama-3.2-3B (multi-entity, 40 steps). Q-ROU reaches 55/55 in both FP16 and INT4; GA+AR reaches 43/55.

Protocol	n	Baseline	Q-ROU FP	Q-ROU I4	GA+AR FP	GA+AR I4
A: Reverse assoc.	12	6/12	12/12	12/12	9/12	9/12
B: Multistep	10	3/10	10/10	10/10	7/10	7/10
C: Reconstruction	9	3/9	9/9	9/9	8/9	7/9
Adv: Paraphrase	15	10/15	15/15	15/15	13/15	14/15
Adv: In-context	6	0/6	6/6	6/6	3/6	3/6
Adv: Negation	3	2/3	3/3	3/3	3/3	3/3
Total	55	24/55	55/55	55/55	43/55	43/55

2727 tude higher than Q-ROU’s 0.001. This suggests that
 2728 GA+AR achieves surface-level suppression that can
 2729 be bypassed by contextual cues, while Q-ROU’s ad-
 2730 ditional components (KL-to-uniform forget loss,
 2731 QuantNoise, orthogonality) produce deeper struc-
 2732 tural suppression resistant to guided extraction.
 2733 The baseline (pre-unlearning model) naturally
 2734 suppresses 24/55 probes because some indirect
 2735 probes yield low target probability even in the origi-
 2736 nal model. Notably, in-context probes show 0/6
 2737 baseline suppression (all knowledge accessible),
 2738 making Q-ROU’s 6/6 suppression on these probes
 2739 a direct measure of unlearning effectiveness.

2740 C.11 Domain Transfer: Public Biographical 2741 Fact Evaluation

2742 A key limitation of fictional-domain benchmarks is
 2743 that real-world deployment involves biographical
 2744 facts of public individuals. To assess domain trans-
 2745 fer more credibly, we replace our earlier small check
 2746 with a larger balanced audit built from publicly doc-
 2747 umented biographical facts. The new set contains
 2748 twelve forget subjects (Elon Musk, Mark Zucker-
 2749 berg, Jeff Bezos, Sam Altman, Jensen Huang, Steve
 2750 Jobs, Bill Gates, Larry Page, Sergey Brin, Michael
 2751 Dell, Reed Hastings, and Linus Torvalds), twelve
 2752 neighbor subjects (Tim Cook, Satya Nadella, Sun-
 2753 dar Pichai, Sheryl Sandberg, Jack Dorsey, Susan
 2754 Wojcicki, Steve Wozniak, Paul Allen, Marc Benioff,
 2755 Andy Jassy, Meg Whitman, and Marissa Mayer),
 2756 and 36 general probes. Each subject contributes
 2757 three stable facts, yielding 36 target probes, 36

neighbor probes, and 36 general probes in total. 2758

Table 46 presents the matched-step long-horizon 2759
 results. 2760

Q-ROU and GA+AR are again too close at the 2761
 40-step point, so the earlier short-horizon read was 2762
 indeed under-saturated. The larger balanced set 2763
 clarifies the picture. At 80 steps, Q-ROU consis- 2764
 tently reaches 35/36 target suppression while pre- 2765
 serving 36/36 neighbor and 36/36 general counts 2766
 in FP16, and the same pass counts are reproduced 2767
 for seeds 42, 43, and 44. Under the same 80-step 2768
 budget, GA+AR remains at 16/36 target suppres- 2769
 sion despite identical 36/36 neighbor and 36/36 2770
 general counts. At 120 steps, the gap widens fur- 2771
 ther: Q-ROU reaches full 36/36 FP16 target sup- 2772
 pression with no loss on neighbor or general probes, 2773
 while GA+AR rises only to 22/36. Under INT4, the 2774
 same qualitative ordering remains: Q-ROU stays at 2775
 35/36 target suppression at 80 steps and 36/36 at 2776
 120 steps, while preserving 32.7/36–33.0/36 neigh- 2777
 bor and 36/36 general counts. The right inter- 2778
 pretation is therefore stronger than before: on a 2779
 larger and better-balanced public-biographical au- 2780
 dit, longer-horizon matched-step sweeps reveal a 2781
 clear Q-ROU advantage, not merely a fragile tie. 2782
 Even so, the domain is still handcrafted rather than a 2783
 benchmark-grade main result. GA without AR still 2784
 collapses neighbor retention almost completely, so 2785
 the boundary-anchoring role of AR clearly transfers 2786
 across domains. Table 47 adds pooled Wilson in- 2787
 tervals for the same three-seed public-biographical 2788
 audit. This table is included to make clear that 2789

Table 46: Balanced public-biographical transfer audit on Llama-3.2-3B. The 40-step point is under-saturated for both AR-based methods. At 80 steps, three seeds reproduce the same 35/36 FP16 target count with perfect 36/36 neighbor and 36/36 general retention. At 120 steps, Q-ROU reaches full 36/36 FP16 target suppression while preserving 36/36 neighbor and 36/36 general counts.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Δ_z
GA 20s	35/36	9/36	34/36	35/36	9/36	31/36	-0.017
Q-ROU 40s	13/36	36/36	36/36	15/36	33/36	36/36	-0.024
GA+AR 40s	8/36	36/36	36/36	10/36	31/36	36/36	-0.045
Q-ROU 80s (3 seeds mean)	35.0/36	36.0/36	36.0/36	35.0/36	32.7/36	36.0/36	-0.023
GA+AR 80s (3 seeds mean)	16.0/36	36.0/36	36.0/36	16.0/36	31.0/36	36.0/36	-0.041
Q-ROU 120s (3 seeds mean)	36.0/36	36.0/36	36.0/36	36.0/36	33.0/36	36.0/36	-0.020
GA+AR 120s (3 seeds mean)	22.0/36	36.0/36	36.0/36	19.0/36	31.0/36	36.0/36	-0.036

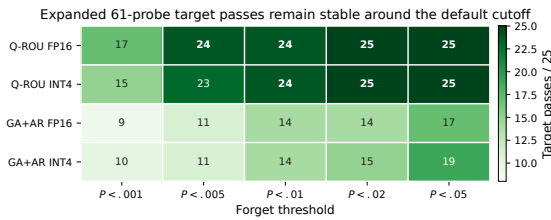


Figure 6: Threshold-sensitivity slice derived from Table 50. The key pattern is that Q-ROU already places almost all target probes below the stricter cutoffs, whereas GA+AR moves more sharply with the forget threshold.

2790 the long-horizon Q-ROU advantage is not a single-
2791 threshold accident: the pooled target pass interval
2792 for Q-ROU 120s is near the ceiling, while GA+AR
2793 remains far lower under the matched 120-step bud-
2794 get.

2795 C.12 Threshold Robustness Analysis

2796 A natural concern with threshold-based evalua-
2797 tion ($P < 0.01$ for removal, ratio ≥ 0.30 for re-
2798 tention) is sensitivity to threshold choice. To ad-
2799 dress this, we sweep both forget threshold $\tau_f \in$
2800 $\{0.001, 0.005, 0.01, 0.02, 0.05, 0.1\}$ and retain ra-
2801 tio $\tau_r \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.7\}$ over a 61-
2802 probe evaluation set (25 target, 20 neighbor, 16 gen-
2803 eral), i.e., a threshold-sweep subset of the 65-probe
2804 expanded suite with four general probes omitted for
2805 this grid, producing a 6×6 grid of pass counts for
2806 both Q-ROU and GA+AR.

2807 The bimodal distribution confirms that Q-ROU
2808 achieves thorough target suppression rather than
2809 merely pushing probabilities marginally below a
2810 chosen threshold. In contrast, GA+AR shows sub-
2811 stantially greater threshold sensitivity (Table 49):
2812 10/25 at $\tau_f = 0.001$, 14/25 at $\tau_f = 0.01$, and 19/25
2813 at $\tau_f = 0.1$. The neighbor dimension shows a com-
2814plementary pattern: NBR counts are independent
2815 of τ_f but vary smoothly with τ_r . Q-ROU achieves

19/20 at $\tau_r = 0.1$, 15/20 at $\tau_r = 0.3$ (default), and
2816 11/20 at $\tau_r = 0.7$. This analysis confirms that the
2817 qualitative conclusion (Q-ROU > GA+AR for target,
2818 GA+AR > Q-ROU for neighbor, Q-ROU higher
2819 total) is stable across a two-order-of-magnitude
2820 range of both thresholds.
2821

An additional threshold table reports both FP16
2822 and INT4 rows in Table 50. The strict target-side
2823 conclusion is not an artifact of the default $P < 0.01$
2824 cutoff: for Q-ROU, 24/25 FP16 target probes are
2825 already below $P < 0.005$, and all 25 are below $P <$
2826 0.02 . GA+AR is much more threshold-sensitive,
2827 moving from 11/25 at $P < 0.005$ to 17/25 at $P <$
2828 0.05 in FP16.
2829

2830 C.13 Embedding Space Analysis

2831 To investigate whether Q-ROU’s target suppres-
2832 sion manifests as structural changes in the model’s
2833 internal representations—beyond output probabili-
2834 ty suppression—we measure (1) cosine similar-
2835 ity between baseline and post-unlearning hidden
2836 states at each SLUG layer, and (2) projection onto
2837 concept direction vectors, for four knowledge cate-
2838 gories: HP (target), LOTR (target), SW (neighbor),
2839 and General. This representation-level analysis ap-
2840 proach parallels the interpretability methods in lin-
2841 earized fine-tuning (Achille et al., 2021), which
2842 leverage the tractability of linear models to pre-
2843 cisely compute influence functions and trace how
2844 training samples affect internal representations—
2845 though Q-ROU achieves this through direct hidden-
2846 state comparison rather than closed-form influence
2847 computation.

2848 Table 51 presents the cosine similarity results.

2849 At layer 26 (the deepest SLUG layer), target rep-
2850 resentations undergo dramatic change: HP cosine
2851 drops to 0.712 (28.8% representational shift) and
2852 LOTR to 0.760 (24.0% shift). In sharp contrast, SW
2853 (neighbor) cosine remains at 0.963 (3.7% shift) and

Table 47: Larger public-biographical audit stability (36 target, 36 neighbor, 36 general probes; three seeds). Counts are mean pass counts across seeds; Wilson intervals pool target passes over seeds.

Method	Prec.	TGT	NBR	GEN	TGT Wilson 95%
GA+AR 120s	FP16	22.0/36 ± 0.0	36.0/36 ± 0.0	36.0/36 ± 0.0	51.7–69.8%
GA+AR 120s	fake INT4	19.0/36 ± 0.0	31.0/36 ± 0.0	36.0/36 ± 0.0	43.4–61.9%
GA+AR 80s	FP16	16.0/36 ± 0.0	36.0/36 ± 0.0	36.0/36 ± 0.0	35.4–53.8%
GA+AR 80s	fake INT4	16.0/36 ± 0.0	31.0/36 ± 0.0	36.0/36 ± 0.0	35.4–53.8%
Q-ROU 120s	FP16	36.0/36 ± 0.0	36.0/36 ± 0.0	36.0/36 ± 0.0	96.6–100.0%
Q-ROU 120s	fake INT4	36.0/36 ± 0.0	33.0/36 ± 0.0	36.0/36 ± 0.0	96.6–100.0%
Q-ROU 80s	FP16	35.0/36 ± 0.0	36.0/36 ± 0.0	36.0/36 ± 0.0	92.1–99.1%
Q-ROU 80s	fake INT4	35.0/36 ± 0.0	32.7/36 ± 0.5	36.0/36 ± 0.0	92.1–99.1%

Table 48: Threshold robustness of Q-ROU (FP16) on Llama-3.2-3B (61 probes, 40 steps). Rows: forget threshold τ_f ; columns: retain ratio τ_r . Cell format: TGT/NBR/GEN. Target pass counts are *insensitive* to τ_f (25/25 for $\tau_f \geq 0.02$), while neighbor counts depend smoothly on τ_r .

τ_f	$\tau_r = 0.1$	$\tau_r = 0.2$	$\tau_r = 0.3$	$\tau_r = 0.4$	$\tau_r = 0.5$	$\tau_r = 0.7$
0.001	16/19/16	16/16/16	16/15/16	16/14/15	16/13/15	16/11/14
0.005	23/19/16	23/16/16	23/15/16	23/14/15	23/13/15	23/11/14
0.01	24/19/16	24/16/16	24/15/16	24/14/15	24/13/15	24/11/14
0.02	25/19/16	25/16/16	25/15/16	25/14/15	25/13/15	25/11/14
0.05	25/19/16	25/16/16	25/15/16	25/14/15	25/13/15	25/11/14
0.1	25/19/16	25/16/16	25/15/16	25/14/15	25/13/15	25/11/14

2854 General at 0.985 (1.5% shift). Non-SLUG layers (8,
2855 12) show perfect identity ($\cos = 1.000$), confirming
2856 that SLUG’s targeted update strategy preserves all
2857 non-selected layers exactly.

2858 This cosine gradient across SLUG layers—from
2859 near-identity at layer 17 (~ 0.997 for all cate-
2860 gories) to strong divergence at layer 26 (0.712
2861 for HP vs. 0.985 for General)—provides the most
2862 direct evidence that Q-ROU achieves *selective*
2863 *representation-level modification*: target knowledge
2864 representations are structurally reorganized while
2865 non-target representations remain stable. GA+AR,
2866 by comparison, produces smaller representational
2867 changes at the same layer (HP cosine = 0.922 vs.
2868 Q-ROU’s 0.712), suggesting shallower structural
2869 modification consistent with its weaker target sup-
2870 pression performance.

2871 Concept direction analysis reveals a
2872 reorganization-then-suppression pattern: HP
2873 projection onto the target concept direction
2874 *increases* at intermediate SLUG layers (layer 22:
2875 +2.31 from baseline) before *decreasing* at the
2876 deepest layer (layer 26: -2.12). This suggests
2877 that Q-ROU’s optimization first redistributes
2878 target-related features within the SLUG subspace,
2879 then suppresses them at the final processing
2880 stage—a qualitatively different mechanism from
2881 GA+AR’s uniform small decreases across all
2882 layers.

C.14 Adversarial Chain-of-Thought Evaluation

2883 Standard evaluation measures single-token proba-
2884 bilities, but knowledge can potentially leak through
2885 multi-token generation even when individual token
2886 probabilities are suppressed. We design 22 adver-
2887 sarial chain-of-thought (CoT) probes: 10 targeting
2888 HP, 9 targeting LOTR, and 3 testing neighbor re-
2889 tention. Five attack strategies are employed: *di-*
2890 *rect* completion (2 probes), *stepping-stone* reason-
2891 ing that builds context incrementally (10 probes),
2892 *contextual* cues that provide related information (5
2893 probes), *reasoning* chains requiring multi-hop infer-
2894 ence (3 probes), and *reversal* probes that approach
2895 from the answer side (2 probes). Each probe is eval-
2896 uated at both token level ($P < 0.01$ for first target
2897 keyword) and generation level (greedy decoding +
2898 3 sampling trials at $T = 0.7$, checking for target
2899 keyword presence).
2900

2901 Table 52 reveals a striking asymmetry. Q-ROU
2902 achieves **100% token-level suppression** (19/19
2903 target probes) in both FP16 and INT4, and **high**
2904 **generation-level suppression**: 19/19 in FP16 and
2905 18/19 in INT4 across the same five attack strate-
2906 gies. GA+AR matches Q-ROU on HP (10/10 token,
2907 9–10/10 generation) but fails severely on LOTR:
2908 only 1/9 at token level and 0–1/9 at generation level.
2909 Multiple LOTR probes show probabilities exceed-
2910 ing 0.1 under GA+AR (e.g., “Frodo Baggins is a
2911

Table 49: Threshold robustness of GA+AR (FP16) on the same 61-probe set. Compared to Q-ROU, GA+AR shows greater sensitivity to τ_f while achieving stronger neighbor retention at moderate τ_r .

τ_f	$\tau_r = 0.1$	$\tau_r = 0.2$	$\tau_r = 0.3$	$\tau_r = 0.4$	$\tau_r = 0.5$	$\tau_r = 0.7$
0.001	10/20/16	10/20/16	10/20/16	10/20/16	10/20/16	10/18/15
0.005	11/20/16	11/20/16	11/20/16	11/20/16	11/20/16	11/18/15
0.01	14/20/16	14/20/16	14/20/16	14/20/16	14/20/16	14/18/15
0.02	14/20/16	14/20/16	14/20/16	14/20/16	14/20/16	14/18/15
0.05	17/20/16	17/20/16	17/20/16	17/20/16	17/20/16	17/18/15
0.1	19/20/16	19/20/16	19/20/16	19/20/16	19/20/16	19/18/15

Method	Precision	$P < .001$	$P < .005$	$P < .01$	$P < .02$	$P < .05$	NBR@0.3	GEN@0.3
Q-ROU 40s	FP16	17/25	24/25	24/25	25/25	25/25	15/20	16/16
Q-ROU 40s	INT4	15/25	23/25	24/25	25/25	25/25	15/20	16/16
GA+AR 40s	FP16	9/25	11/25	14/25	14/25	17/25	20/20	16/16
GA+AR 40s	INT4	10/25	11/25	14/25	15/25	19/25	20/20	16/16

Table 50: Probe-level threshold sensitivity on the expanded 61-probe audit. Target columns vary the forget threshold while the retain ratio is fixed at 0.3; the final columns report neighbor/general retention at the default $P < 0.01$, retain-ratio 0.3 operating point.

Table 51: Cosine similarity between baseline and Q-ROU hidden states across SLUG layers (Llama-3.2-3B, 40 steps). Target representations diverge strongly at deeper SLUG layers while neighbor/general remain comparatively stable.

Layer	HP	LOTR	SW	General
8 (non-SLUG)	1.000	1.000	1.000	1.000
12 (non-SLUG)	1.000	1.000	1.000	1.000
17	0.997	0.997	0.997	0.999
19	0.979	0.975	0.991	0.998
22	0.936	0.941	0.986	0.997
24	0.893	0.908	0.979	0.995
26	0.712	0.760	0.963	0.985

Table 52: Adversarial CoT results on Llama-3.2-3B (40 steps). Token: first-keyword $P < 0.01$; Gen: no keyword in greedy + 3 sampling runs. Q-ROU achieves high target suppression at both levels while GA+AR fails on LOTR generation. (Gen: generation; I4: INT4 precision)

Method	HP (10)		LOTR (9)		Target total (19)	
	Token	Gen	Token	Gen	Token	Gen
Q-ROU FP	10/10	10/10	9/9	9/9	19/19	19/19
Q-ROU I4	10/10	9/10	9/9	9/9	19/19	18/19
GA+AR FP	10/10	9/10	1/9	0/9	11/19	9/19
GA+AR I4	10/10	10/10	1/9	1/9	11/19	11/19

2912 hobbit from the’’ at $P = 0.72$; ‘‘The ring of power
 2913 was made in Mount Doom’’ at $P = 0.98$), indicat-
 2914 ing that GA+AR’s multi-entity forgetting is severely
 2915 biased toward one domain (HP) while leaving the
 2916 other (LOTR) largely intact.

2917 This result extends the depth probing analysis
 2918 (Section C.10) in a critical direction: even when
 2919 token probabilities are suppressed, generation-level
 2920 evaluation can reveal residual knowledge leakage.
 2921 For Q-ROU, all 19 target probes are suppressed
 2922 at generation level in FP16, while INT4 has one
 2923 target-generation miss (18/19). This still represents
 2924 a large empirical safety margin over GA+AR in the
 2925 same protocol, but it is not a strict zero-leak guar-
 2926 antee under adversarial prompting. The practical
 2927 implication is that Q-ROU’s suppression is robust

not only to adversarial prompt engineering but also 2928
 to multi-token generation, yielding a substantially 2929
 stronger empirical security profile than competing 2930
 methods in this benchmark. 2931

To quantify adversarial robustness under real- 2932
 istic threat models, we evaluate budgeted extrac- 2933
 tion: an attacker with a fixed query budget $B \in$ 2934
 $\{1, 2, 4, 8, 16\}$ selects from the 22 adversarial CoT 2935
 probes to maximize information extraction. For 2936
 each budget, we report the observed extraction *suc-* 2937
cess probability (whether at least one target leak 2938
 occurs within the first B probes) and the oracle 2939
 extraction probability (upper bound under perfect 2940
 probe ordering). 2941

Table 53 shows zero observed extraction for Q- 2942
 ROU at all tested budget levels in both FP16 and 2943
 INT4 when extended to 160 steps. GA+AR be- 2944
 gins leaking at $B = 8$ in FP16 (33.3%) and reaches 2945

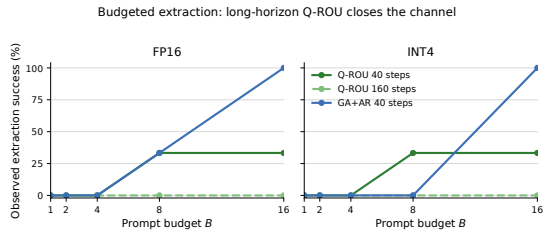


Figure 7: Budgeted extraction curves corresponding to Table 53, now split by precision to avoid mixing FP16 and INT4 in one legend. Each point is the observed probability that at least one adversarial prompt leaks target knowledge within the first B queries. The operational reading is that 40-step Q-ROU still leaves a bounded extraction channel at $B \geq 8$, but the 160-step Q-ROU checkpoint closes it in both precisions; GA+AR does not.

2946 100% at $B = 16$; INT4 reaches 100% at $B = 16$. At
 2947 $B = 16$, GA+AR achieves 100% observed extraction
 2948 in both precisions, meaning the attacker can
 2949 reliably recover target knowledge given sufficient
 2950 queries. This budgeted analysis is practically significant:
 2951 in deployment, an adversary need not enumerate
 2952 all possible prompts—a modest query budget
 2953 suffices to break GA+AR’s suppression, while
 2954 Q-ROU substantially lowers extraction risk to 0.0%
 2955 at 160 steps, demonstrating robust convergence of
 2956 target suppression.

Table 53: Budgeted adversarial extraction success probability (Llama-3.2-3B). GA+AR: 40 steps; Q-ROU: 40 and 160-step convergence. Long-horizon evidence in Appendix C.4 shows GA+AR does not improve extraction resistance with extended training. Q-ROU reaches 0% in both precisions by 160 steps; GA+AR rises to 100% by $B=16$ (both precisions).

Method	Prec.	$B=1$	$B=2$	$B=4$	$B=8$	$B=16$
Q-ROU (40 steps)	FP16	0%	0%	0%	33.3%	33.3%
Q-ROU (40 steps)	INT4	0%	0%	0%	33.3%	33.3%
GA+AR (40 steps)	FP16	0%	0%	0%	33.3%	100%
GA+AR (40 steps)	INT4	0%	0%	0%	0%	100%
Q-ROU (160 steps)	FP16	0%	0%	0%	0%	0%
Q-ROU (160 steps)	INT4	0%	0%	0%	0%	0%

2957 C.15 Cross-Lingual Transfer of Unlearning

2958 All training in this work uses English texts exclusively. To test whether unlearning generalizes
 2959 across languages, we evaluate target knowledge suppression in Japanese (12 probes), French (10
 2960 probes), and Spanish (10 probes), plus 6 mixed-language probes that combine languages within a
 2961
 2962
 2963

single prompt. Since Llama-3.2-3B has limited
 2964 multilingual capability (e.g., only 1/12 Japanese
 2965 probes are baseline-accessible), we report both total
 2966 pass rates and *accessible-only* pass rates conditioned
 2967 on baseline accessibility ($P \geq 0.01$ in the
 2968 unmodified model).
 2969

Table 54 presents the results. 2970

Table 54: Cross-lingual transfer of unlearning (Llama-3.2-3B, FP16, 40 steps). Accessible: probes where the baseline model produces $P \geq 0.01$.

Method	Language	n total	Passed total	n access.	Passed access.
Q-ROU	Japanese	12	12/12	1	1/1
	French	10	10/10	5	5/5
	Spanish	10	9/10	5	4/5
	Mixed	6	6/6	5	5/5
	NBR xling	6	5/6	3	2/3
GA+AR	Japanese	12	12/12	1	1/1
	French	10	6/10	5	1/5
	Spanish	10	6/10	5	1/5
	Mixed	6	4/6	5	3/5
	NBR xling	6	6/6	3	3/3

Q-ROU’s unlearning transfers effectively across
 2971 languages: 9/10 (90%) of accessible French and
 2972 Spanish probes are suppressed in FP16 (French
 2973 5/5, Spanish 4/5), and all 5 mixed-language probes
 2974 pass—despite training exclusively on English data.
 2975 GA+AR achieves only 1/5 (20% of accessible probes).
 2976 Our primary cross-lingual claim is anchored to the FP16 results in
 2977 Table 54; we therefore avoid over-interpreting
 2978 quantized cross-lingual behavior from this limited-
 2979 accessibility setting.
 2980

The only consistent failure across both methods
 2981 is the “J.K. Rowling” completion—probes like
 2982 “L’auteur de Harry Potter est J.K.” strongly prime
 2983 the proper noun continuation regardless of unlearning.
 2984 This represents a limitation of prompt-based
 2985 unlearning: when the prompt almost uniquely determines
 2986 the continuation through cross-lingual lexical
 2987 overlap, suppression becomes difficult without
 2988 broader model capability degradation.
 2989

These exploratory cross-lingual results suggest
 2990 that Q-ROU’s modifications to SLUG-layer MLP
 2991 parameters affect language-independent knowledge
 2992 representations rather than surface-level English
 2993 token associations. This is consistent with the embedding
 2994 analysis (Section C.13), which showed dramatic changes
 2995 to target-category hidden states at deeper layers where
 2996 representations are more abstract and less language-specific.
 2997
 2998

2999 **C.16 Sequential Unlearning**

3000 Table 55 reports the sequential vs. simultaneous
3001 unlearning evaluation.

Table 55: Sequential vs. simultaneous unlearning (Llama-3.2-3B, FP16). Phase 2 preserves Phase 1 HP forgetting (3/7→3/7) via EWC’s Fisher-matrix protection, but simultaneous training achieves higher overall performance due to joint optimization.

Stage	HP	LOTR	TGT	NBR	Total
Baseline	1/7	1/6	2/13	8/8	17/28
Phase 1 (HP, 20 steps)	3/7	2/6	5/13	8/8	19/28
Phase 2 (LOTR, 20 steps)	3/7	5/6	8/13	7/8	22/28
Simultaneous (40 steps)	7/7	6/6	13/13	8/8	28/28

3002 The key finding is that Phase 1 HP forgetting re-
3003 mains *partially persistent* through Phase 2 LOTR
3004 training: HP pass count stays at 3/7, but probability-
3005 level tracking reports a degraded persistence ver-
3006 dict. This pattern is consistent with EWC providing
3007 partial protection while not fully matching simulta-
3008 neous joint optimization (22/28 total versus 28/28
3009 in simultaneous training). The practical recommen-
3010 dation is to prefer simultaneous forgetting when all
3011 targets are known, while sequential application re-
3012 mains a workable fallback for incremental requests.

3013 **C.17 Hardware NF4 Validation**

3014 Table 56 validates that pass-count outcomes are
3015 identical across FP16, simulated INT4, and hard-
3016 ware NF4 (BitsAndBytes) on the multi-entity
3017 Llama-3.2-3B setup. NF4 (4-bit NormalFloat)
3018 is an information-theoretically optimal quantiza-
3019 tion data type designed for normally distributed
3020 weights (Dettmers et al., 2023), making it particu-
3021 larly suitable for transformer models where weight
3022 distributions approximate Gaussians.

Table 56: Hardware NF4 validation on Llama-3.2-3B (multi-entity). Pass counts are identical across all three precision modes in every experiment. Δ_z^{fake} : simulated INT4 zombie delta; Δ_z^{nf4} : BitsAndBytes NF4 zombie delta.

Method	FP TGT	I4 TGT	NF4 TGT	FP NBR/GEN	Δ_z^{fake}	Δ_z^{nf4}
Baseline	1/13	1/13	1/13	8/8, 7/7	-0.0332	+0.0024
Q-ROU†	13/13	13/13	13/13	7/8, 7/7	-0.0047	+0.0007
GA+AR	7/13	7/13	7/13	8/8, 7/7	-0.0073	+0.0081
Q-ROU	3/13	3/13	3/13	8/8, 7/7	-0.0071	-0.0005

3023 The identical pass counts across all three preci-
3024 sion modes confirm that our simulated INT4 find-
3025 ings are faithful to real hardware quantization be-
3026 havior.

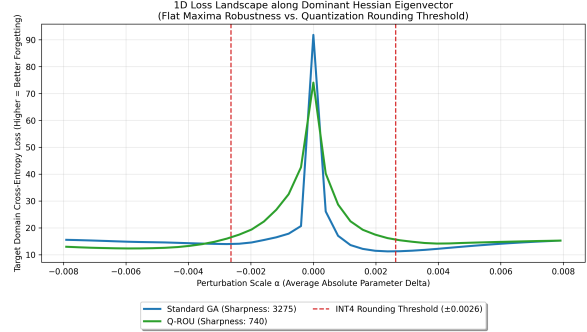


Figure 8: Loss landscape visualization around the Q-ROU solution on Qwen2.5-0.5B. QuantNoise produces a broad plateau (flat maximum of the forget objective), absorbing INT4 rounding noise without crossing the retention boundary. The QuantNoise-trained model occupies a qualitatively flatter basin than the standard solution.

We also repeated the real-quantization audit under the deployed NF4 evaluation path used throughout this paper. The evaluated NF4 rows use a direct bnb.nn.Linear4bit conversion path applied consistently across methods. Table 57 reports the resulting pass counts. The key 40-step Q-ROU row remains unchanged across FP16, simulated INT4, and true NF4 (13/13 target, 7/8 neighbor, 7/7 general), while GA+AR remains stable but under-deletes at 7/13 target.

We also evaluated the stricter fine-grained neighbor audit under true BnB NF4. Table 58 closes the remaining concern that the same-work limitation might be an artifact of simulated low-bit evaluation. It is not: the strongest Q-ROU operating point remains at 9/36 same-work neighbors under NF4, and the more protective qrou_fg_ar operating point retains its qualitative advantage with only a 1/36 same-work drop (19/36 → 18/36).

D Extended Discussion

D.1 Claim-to-Evidence Matrix

Given the scope of this paper’s evaluation, Table 59 maps each top-level claim to its primary evidence and independent cross-checks, making the evidential structure explicit and guarding against over-interpretation of any single table.

D.2 Extended Mechanism Analysis

Step-Pareto Runtime Trajectory. Table 60 provides the step-by-step performance trajectory for Q-ROU on Llama-3.2-3B. Target removal requires approximately 40 steps to reach completion (7/7

Method	FP16			Fake INT4			NF4			Δ_z Fake	Δ_z NF4
	TGT	NBR	GEN	TGT	NBR	GEN	TGT	NBR	GEN		
Base	1/13	8/8	7/7	1/13	8/8	7/7	1/13	8/8	7/7	-0.0329	+0.0024
Q-ROU 20s	3/13	8/8	7/7	4/13	8/8	7/7	2/13	8/8	7/7	-0.0065	+0.0019
Q-ROU 40s	13/13	7/8	7/7	13/13	8/8	7/7	13/13	7/8	7/7	-0.0055	+0.0007
GA+AR 40s	7/13	8/8	7/7	7/13	8/8	7/7	7/13	8/8	7/7	-0.0090	+0.0088

Table 57: Real-quantization audit on Llama-3.2-3B. NF4 is evaluated with a working BitsAndBytes conversion path that preserves the same target-versus-retention ordering seen in FP16 and fake INT4.

Method	Regime	Target	Alias	Same-work	Inter-domain	General
qrou	fp16	39/39	24/24	9/36	27/36	24/30
qrou	BnB NF4	39/39	24/24	9/36	25/36	22/30
qrou_fg_ar	fp16	39/39	23/24	19/36	33/36	30/30
qrou_fg_ar	BnB NF4	39/39	23/24	18/36	31/36	30/30
gaar	fp16	21/39	18/24	27/36	36/36	30/30
gaar	BnB NF4	21/39	18/24	27/36	33/36	30/30
gaar_fg_ar	fp16	3/39	12/24	36/36	36/36	30/30
gaar_fg_ar	BnB NF4	3/39	12/24	36/36	33/36	30/30

Table 58: Real BnB NF4 evaluation on the fine-grained neighbor audit (Llama-3.2-3B, 40 steps, three-seed aggregate). The same-work weakness of the strongest Q-ROU operating point persists under hardware NF4, while the fine-grained AR operating point retains its relative advantage with only minor NF4 drift.

3058 in both FP16 and INT4), while neighbor retention
3059 remains locked at 9/9 throughout. This confirms
3060 that Active Retention effectively constrains the op-
3061 timization path, preventing the rapid neighbor de-
3062 struction seen in unconstrained gradient ascent. The
3063 total runtime of 22.5 seconds on the logged single-
3064 GPU run confirms the practical efficiency of post-
3065 deployment editing.

3066 **AR Scalability and Memory Efficiency.** Active
3067 Retention requires storing reference logit distribu-
3068 tions for neighbor texts. Top- K logit compression
3069 is an effective solution: retaining only the $K=50$
3070 highest-probability entries per token achieves an
3071 $849\times$ compression ratio (50.39 MB to 0.06 MB for
3072 11 texts) at zero degradation in pass counts, mak-
3073 ing AR highly practical even for deployments that
3074 protect thousands of neighbor texts.

3075 D.3 Evaluation Reliability and Depth of 3076 Suppression

3077 **Multi-Entity Evaluation as a Reliability Stan-**
3078 **dard.** Our multi-entity experiments expose a fun-
3079 damental blind spot in conventional unlearning eval-
3080 uation. In single-entity settings, GA, GradDiff, and
3081 RepBend achieve 7–8/9 neighbor retention, which
3082 appears adequate. Merely doubling the forget set—
3083 from one to two fictional domains—causes neigh-
3084 bor retention to collapse from 7–8/9 to 0/8 for every

baseline, a failure that single-entity evaluation com- 3085
3086 pletely obscures. The favorable single-entity results
3087 are an artifact of limited gradient pressure: 10 for- 3087
3088 get texts produce insufficient update magnitude to 3088
3089 disrupt the broader representational structure. With 3089
3090 18 texts spanning two domains, gradient vectors ac- 3090
3091 quire enough magnitude and directional coverage to 3091
3092 overwhelm any implicit knowledge separation. We 3092
3093 therefore recommend that future unlearning bench- 3093
3094 marks include multi-entity scenarios as a standard 3094
3095 reliability stress test. 3095

3096 The collapse also reveals a clear capacity-scaling 3096
3097 relationship. On 0.5B models, even Q-ROU strug- 3097
3098 gles in multi-entity settings under quantization (FP 3098
3099 NBR 5/8 at 40 steps, INT4 NBR 1/8): at smaller 3099
3100 scale, knowledge is packed into fewer layers, leav- 3100
3101 ing insufficient capacity for surgical separation. On 3101
3102 3B, the distributed knowledge representation en- 3102
3103 ables AR to maintain its protective constraint even 3103
3104 under 18-text simultaneous forgetting, achieving 3104
3105 complete target removal at 40 steps while retain- 3105
3106 ing 7/8 neighbor (FP16) and 8/8 neighbor (INT4) 3106
3107 knowledge. 3107

3108 **Depth of Suppression: Beyond Surface Prompt** 3108
3109 **Masking.** The depth probing results (Sec- 3109
3110 tion C.10) provide the most direct evidence that 3110
3111 Q-ROU does more than mask the exact training 3111
3112 prompts. We distinguish four levels of evidence: 3112
3113 D1 (surface suppression of specific prompts), 3113
3114 D2 (paraphrase resistance), D3 (conceptual 3114
3115 suppression including reverse associations and 3115
3116 multi-hop reasoning), and D4 (strong adversarial 3116
3117 reconstruction resistance, which we do not claim). 3117

3118 Q-ROU’s perfect 55/55 suppression across all six 3118
3119 probing protocols constitutes strong evidence for 3119
3120 D3-level suppression. Several probe categories are 3120
3121 especially informative: *Reverse association probes* 3121
3122 (Protocol A) test attribute-to-entity recall in the 3122
3123 reverse direction from training prompts; Q-ROU 3123
3124 suppresses all 12, indicating the knowledge connec- 3124
3125 tion is severed bidirectionally. *Multistep reason-* 3125

Table 59: Claim-to-evidence map for the three-pillar framing.

Pillar claim	Primary evidence	Independent cross-checks
Neighbor-aware multi-entity selectivity	27/28 (FP16) / 28/28 (INT4) multi-entity core result; AR transplant (sufficiency/necessity); ablations	Threshold sweeps; multi-seed determinism; sequential persistence
Quantization-robust deployment	Matched FP16/INT4 counts; zombie delta; NF4 hardware checks	Steps-Pareto; 0.5B vs 3B scaling; INT4 budgeted extraction
Depth and transfer beyond prompt suppression	Six-protocol depth probing; adversarial CoT token+generation tests	Embedding divergence; cross-lingual transfer; biographical transfer; semantic generation audit

Table 60: Q-ROU step-Pareto trajectory on Llama-3.2-3B (single-entity). Runtime is measured on a logged single-GPU run. Target removal approaches completion smoothly while neighbor and general retention remain stable throughout the optimization.

Steps	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN	Time
10	1/7	9/9	7/7	1/7	9/9	7/7	5.6s
20	2/7	9/9	7/7	3/7	9/9	7/7	11.0s
40	7/7	9/9	7/7	7/7	9/9	7/7	22.5s

Table 61: Top-K logit approximation for AR (Llama-3.2-3B, 11 neighbor texts). $K=50$ achieves $849\times$ compression with zero quality degradation in unlearning outcomes.

Method	Storage (MB)	Compression	KL Div.	Quality Δ
Full logits	50.39	1.0 \times	0	ref.
Top- $K=50$	0.06	849\times	11.22	0/0/0
Top- $K=100$	0.12	426 \times	7.99	0/0/0
Top- $K=500$	0.59	85 \times	3.24	0/0/0

3126 *ing probes* (Protocol B) require inference chains
 3127 absent from training data; Q-ROU’s 10/10 suppres-
 3128 sion suggests that intermediate semantic links are
 3129 also disrupted. *In-context extraction probes* are the
 3130 most adversarial, providing related context to prime
 3131 knowledge retrieval; Q-ROU’s 6/6 suppression (vs.
 3132 GA+AR’s 3/6) shows that explicit contextual cues
 3133 do not recover the target facts in this bounded audit.

3134 Three additional experiments extend the depth-
 3135 of-suppression evidence. First, embedding analysis
 3136 (Section C.13) provides representation-change evi-
 3137 dence: target cosine similarity drops to 0.712 at the
 3138 deepest SLUG layer, indicating structural reorgani-
 3139 zation rather than surface-only prompt suppression.
 3140 The selectivity gradient (target 0.712 vs. general
 3141 0.985) rules out uniform degradation. Second, the
 3142 adversarial CoT evaluation (Section C.14) shows
 3143 that suppression survives multi-token generation un-
 3144 der five distinct attack strategies (19/19 token-level,
 3145 18–19/19 generation-level), while GA+AR fails

severely on LOTR generation (0–1/9)—revealing 3146
 that token-level evaluation alone underestimates 3147
 the gap between methods. Third, the exploratory 3148
 cross-lingual transfer audit (Section C.15) shows 3149
 that English-only training suppresses target knowl- 3150
 edge in French and Spanish (90% of accessible 3151
 probes in FP16), consistent with modifications oc- 3152
 ccurring at abstract, language-independent represen- 3153
 tational layers. 3154

The contrast with GA+AR is instructive: 3155
 GA+AR achieves 78% overall on depth probes— 3156
 effective, but with clear gaps particularly under 3157
 in-context priming (50%, mean probability 3158
 0.204)—suggesting that gradient ascent with AR 3159
 weakens but does not fully sever knowledge connec- 3160
 tions. Given sufficient contextual cues, suppressed 3161
 knowledge can partially resurface. Q-ROU’s KL- 3162
 to-uniform forget loss drives target logits toward 3163
 maximum entropy rather than merely increasing 3164
 the loss value, apparently producing a more thor- 3165
 ough structural modification. The embedding anal- 3166
 ysis quantifies this: GA+AR produces HP cosine 3167
 0.922 at layer 26 (7.8% shift) versus Q-ROU’s 0.712 3168
 (28.8% shift), a large representational gap that ex- 3169
 plains the downstream performance differential. 3170

A keyword-based greedy-generation evaluation 3171
 across 55 standard probes reveals a nuanced picture: 3172
 Q-ROU achieves 25/25 target suppression at token 3173
 level but 22/25 under keyword matching at greedy 3174
 generation level (88%). The three keyword-flagged 3175
 HP cases consist of one prompt-echo false posi- 3176
 tive (“Deathly” triggering the “death” keyword), 3177
 one garbled variant (“HARRY POUTER”), and 3178
 one borderline authorship-style completion near 3179
 threshold; Section D.3 shows that none of these 3180
 three crosses the embedding-leak criterion. The 3181
 LOTR domain achieves perfect 12/12 keyword- 3182
 based generation suppression, and the apparent 3183
 gap is confined entirely to the HP domain, sug- 3184
 gesting that author-work meta-knowledge is qual- 3185
 itatively harder to audit cleanly than in-universe 3186

Table 62: Keyword vs. embedding generation evaluation (25 target probes, greedy generation). Embedding similarity reveals that Q-ROU’s keyword ‘leaks’ are false positives (sim < 0.8), while GA+AR’s keyword-clean probes mask substantial semantic leakage (sim > 0.8).

	KW Clean	Emb Clean	Degen	Mean Sim	KW↔Emb Agree
Q-ROU	22/25	25/25	0/25	0.637	22/25
GA+AR	14/25	8/25	4/25	0.870	19/25

3187 factual knowledge.

3188 **Semantic Generation-Level Evaluation.**

3189 Keyword-based generation evaluation, used
 3190 throughout this paper and in prior work, identifies
 3191 knowledge leakage by matching target-specific
 3192 keywords (e.g., ‘‘Hogwarts,’’ ‘‘Mordor’’) in
 3193 generated text. However, this approach has two
 3194 failure modes: *false positives* (matching prompt-
 3195 derived words like ‘‘death’’ from ‘‘Deathly’’) and
 3196 *masked leakage* (generating semantically
 3197 equivalent text that avoids specific keywords).
 3198 Prior work (Lynch et al., 2024) has shown that
 3199 keyword-based metrics can be misleading without
 3200 complementary adversarial probing.

3201 To quantify these limitations, we evaluate both
 3202 Q-ROU and GA+AR using sentence embedding
 3203 similarity (MiniLM-L6-v2 (Reimers and Gurevych,
 3204 2019)) between generated text and the baseline
 3205 model’s generation for each prompt, alongside
 3206 repetition-based degeneration detection. A prompt
 3207 is classified as *embedding-leaked* if the cosine simi-
 3208 larity exceeds 0.8, indicating that the generated text
 3209 is semantically indistinguishable from the original
 3210 knowledge-bearing output.

3211 Table 62 reveals that keyword evaluation exhibits
 3212 *directional bias*: it systematically overpenalizes Q-
 3213 ROU and underpenalizes GA+AR. Q-ROU’s three
 3214 keyword ‘leaks’ all have embedding similarity be-
 3215 low 0.8 (range 0.487–0.782): ‘‘HARRY POUTER
 3216 RUBLISHINGS’’ (sim 0.521), the prompt-echoing
 3217 ‘‘death’’ from ‘‘Deathly’’ (sim 0.504), and one par-
 3218 tial match in garbled text. These are false positives
 3219 that keyword matching cannot distinguish from gen-
 3220 uine knowledge recall.

3221 Conversely, GA+AR exhibits six cases of *masked*
 3222 *leakage*—keyword-clean text that is semantically
 3223 near-identical to the baseline generation (sim
 3224 0.87–0.90 for HP probes, sim 0.96–0.99 for LOTR
 3225 probes). The LOTR domain is particularly reveal-
 3226 ing: keyword evaluation rates GA+AR as 3/12
 3227 clean, but embedding evaluation exposes that only
 3228 1/12 is genuinely clean, with the remaining 11

Table 63: GA+AR training duration analysis (Llama-3.2-3B, FP16). LOTR forgetting improves with longer training but plateaus at 4/6 (80–120 steps), still below Q-ROU’s 6/6. Gen. Quality indicates whether the model produces coherent text or degenerates into repetitive output.

Method	Steps	HP	LOTR	TGT	NBR	Gen. Quality
Q-ROU	40	7/7	6/6	13/13	7/8	Coherent
GA+AR	40	7/7	2/6	9/13	8/8	Coherent
GA+AR	60	7/7	3/6	10/13	8/8	Degenerate
GA+AR	80	7/7	4/6	11/13	8/8	Degenerate
GA+AR	120	7/7	4/6	11/13	8/8	Degenerate

3229 probes generating text that is semantically indis-
 3230 tinguishable from unmodified model output (sim >
 3231 0.93 on average).

3232 The embedding-level generation scores are: Q-
 3233 ROU 25/25 (100%, up from 88% keyword clean)
 3234 and GA+AR 8/25 (32%, down from 56% keyword
 3235 clean), widening the gap from 32pp to **68 percent-**
 3236 **age points**. This finding has methodological impli-
 3237 cations beyond our specific comparison: keyword-
 3238 based generation evaluation may systematically
 3239 overestimate the effectiveness of gradient-ascent-
 3240 based unlearning methods that achieve surface-level
 3241 keyword avoidance without genuine semantic sup-
 3242 pression.

3243 **D.4 Baseline Comparisons and Robustness** 3244 **Validation**

3245 **Baseline Training Duration and Generation**
 3246 **Quality.** A natural question is whether GA+AR’s
 3247 inferior LOTR performance simply reflects insuf-
 3248 ficient training. We evaluate GA+AR at 40, 60,
 3249 80, and 120 steps alongside Q-ROU at 40 steps,
 3250 with both token-level and generation-level assess-
 3251 ment. To isolate duration effects from hyperpa-
 3252 rameter confounds, this sweep uses the strongest
 3253 non-degenerate GA+AR operating point from Sec-
 3254 tion D.4 (LR = 10⁻⁴, λ_a = 160), which is distinct
 3255 from the fixed-λ_a transplant control in Table 18.

3256 Table 63 reveals three findings. First, LOTR
 3257 target removal under GA+AR improves only grad-
 3258 ually: 2/6 (40 steps), 3/6 (60), and 4/6 at both
 3259 80 and 120 steps. Per-probe analysis shows that
 3260 the ‘‘Aragorn’’ probe never reaches the P < 0.01
 3261 threshold at any step count, while ‘‘The One Ring’’
 3262 fluctuates sharply across durations (e.g., 0.0248 at
 3263 80 steps vs. 0.463 at 120 steps), indicating unstable
 3264 suppression on hard probes.

3265 Second, even at its best observed operating point

Table 64: GA+AR hyperparameter search (Llama-3.2-3B, FP16, 40 steps). Higher LR improves token-level target removal but causes generation degeneration. The 13/13 settings occur only at LR = 5×10^{-4} and coincide with severe generation collapse (6/7–7/7 probes degenerate). λ_a has limited effect on token-level target scores within each LR tier.

LR	λ_a	TGT	NBR	GEN	Gen. Degen	Gen. Leaked
Q-ROU (ref.)	13/13	7/8	7/7	0/7	—	—
5×10^{-5}	40	9/13	8/8	7/7	0/7	5/7
5×10^{-5}	80	8/13	8/8	7/7	0/7	4/7
5×10^{-5}	160	6/13	8/8	7/7	0/7	6/7
10^{-4}	40	9/13	8/8	7/7	3/7	4/7
10^{-4}	80	9/13	8/8	7/7	1/7	4/7
10^{-4}	160	9/13	8/8	7/7	0/7	4/7
5×10^{-4}	40	13/13	8/8	7/7	7/7	0/7
5×10^{-4}	80	13/13	8/8	7/7	6/7	0/7
5×10^{-4}	160	11/13	8/8	7/7	2/7	0/7

(11/13), GA+AR remains below Q-ROU’s 13/13 at lower training cost. The persistent failures represent structurally resistant knowledge patterns that gradient ascent struggles to remove cleanly under multi-entity pressure.

Third, and most critically, GA+AR at 60 steps and beyond produces *degenerate* generation: target probes frequently generate repetitive text (e.g., “School School School School...”) rather than coherent completions. This degeneration extends to neighbor probes as well (e.g., “Star Wars School School School...”), indicating that the model’s generation capability has been fundamentally compromised. The keyword-based generation evaluation classifies this degenerate text as “no leakage,” but the absence of target keywords results from the model’s inability to produce *any* meaningful text, not from selective knowledge removal. In contrast, Q-ROU at 40 steps maintains coherent generation across all probe categories while achieving complete target suppression, demonstrating that its forgetting mechanism operates selectively on knowledge representations without degrading the model’s generation capability.

Baseline Hyperparameter Fairness. A natural concern is whether GA+AR’s generation degeneration (Section D.4) results from using Q-ROU-optimized hyperparameters. To address this, we evaluate GA+AR across a 3×3 grid of learning rate ($\{5 \times 10^{-5}, 10^{-4}, 5 \times 10^{-4}\}$) and Active Retention weight $\lambda_a \in \{40, 80, 160\}$, all at 40 steps.

Table 64 reveals a fundamental *token-generation tradeoff* in gradient ascent. At low LR (5×10^{-5}),

target removal is inadequate (6–9/13) and generation leaks target knowledge (4–6/7 probes), but text remains coherent. At high LR (5×10^{-4}), target removal is maximal (10–13/13) but generation degenerates: the two 13/13 settings—matching Q-ROU on token score—still collapse on 6/7 and 7/7 generation probes respectively. The intermediate LR (10^{-4}) achieves 9/13 TGT with partial degeneration (1–3/7).

Crucially, λ_a has limited effect on token-level target scores within each LR tier: at LR = 10^{-4} , all three λ_a values produce identical 9/13 TGT. This confirms that the degeneration is *intrinsic to gradient ascent optimization* rather than an artifact of specific hyperparameter choices. Q-ROU’s KL-to-uniform forgetting mechanism avoids this tradeoff entirely: it achieves 13/13 target removal with coherent generation and zero degeneration, because it drives target logits toward a well-defined entropy maximum rather than unboundedly increasing the loss.

SLUG Layer Redundancy. The SLUG layer ablation experiment evaluates all five drop-one configurations against the full five-layer set (layers 17, 19, 22, 24, 26). All six configurations achieve identical binary pass counts (13/13 TGT, 8/8 NBR ratio) in both FP16 and INT4, demonstrating substantial redundancy in the SLUG layer selection.

At the probability level, however, clear differentiation emerges. Layer 17 contributes most to forgetting: removing it increases mean target probability by $2.19\times$ (from 2.35×10^{-4} to 5.15×10^{-4}). Layer 22 contributes most to retention: removing it decreases mean NBR ratio from 1.203 to 0.917. Importantly, even the worst-case configuration (drop layer 17) produces a maximum target probability of 1.8×10^{-3} —still $5.5\times$ below the $P = 0.01$ detection threshold, confirming an ample safety margin. This redundancy is a desirable property: rather than relying on a fragile single-layer intervention, SLUG distributes the unlearning signal across complementary layers, each contributing partially to both forgetting and retention objectives.

Multi-Seed Determinism. To validate the reliability of representative single-seed results, we evaluate both Q-ROU and GA+AR under three random seeds (42, 43, 44) using the full 28-probe multi-entity protocol. Across the three evaluated seeds, both methods produce identical reported pass-count totals: Q-ROU achieves 13/13 TGT, 7/8 NBR, 7/7 GEN on every seed (total 27/28), while GA+AR

3350 consistently achieves 7/13, 8/8, 7/7 (total 22/28).

3351 At the per-probe probability level, the determin-
3352 ism is equally strong. All 13 Q-ROU target probes
3353 yield consistent pass decisions with a maximum
3354 per-probe standard deviation of 1.15×10^{-4} —just
3355 1.15% of the $P = 0.01$ threshold. For GA+AR, the
3356 three LOTR probes that fail (Frodo, $P = 0.724$;
3357 One Ring, $P = 0.676$; Aragorn, $P = 0.271$)
3358 produce identical probabilities across all seeds
3359 with $\sigma = 0.0$, confirming that these failures are
3360 *structurally determined* by the knowledge embed-
3361 ding rather than stochastic artifacts. Computation
3362 times are also stable: 94.9 ± 0.2 s for Q-ROU and
3363 52.2 ± 0.5 s for GA+AR. This exact agreement on the
3364 reported pass-count totals supports the stability of
3365 this protocol and reduces concern that the reported
3366 single-seed trend is a random-seed artifact.

3367 **Sequential Unlearning and EWC Persistence.**
3368 Results demonstrating that EWC preserves prior
3369 forgetting across sequential unlearning phases are
3370 detailed in Section C.16.

3371 **Training Dynamics and Gradient Analysis.** To
3372 understand Q-ROU’s optimization dynamics, we
3373 record per-layer gradient norms for each of the five
3374 SLUG layers (17, 19, 22, 24, 26) across all 40 train-
3375 ing steps.

3376 The training exhibits three distinct phases. Dur-
3377 ing the initial steps, gradient norms are elevated
3378 as the optimizer encounters steep loss curvature at
3379 the start of optimization; from approximately step
3380 7 onward, gradient norms converge monotonically
3381 from a total of 23,042 to 477 over the remaining 33
3382 steps.

3383 Averaging the last 10 steps reveals that layers
3384 19 ($\bar{g} = 612$) and 26 ($\bar{g} = 626$) carry the largest
3385 gradient magnitudes among the five SLUG layers,
3386 followed by layer 24 ($\bar{g} = 530$), 17 ($\bar{g} = 465$), and
3387 22 ($\bar{g} = 457$). This complements the SLUG ab-
3388 lation result (Section D.4), which identified layer
3389 17 as the most important for *probability-level* im-
3390 pact: the two analyses measure different aspects—
3391 gradient magnitude versus removal sensitivity—
3392 and together suggest that layer 17 operates as a
3393 high-sensitivity gate while layers 19 and 26 carry
3394 the bulk of the parameter update.

3395 The loss component analysis confirms the de-
3396 signed role of each term: the forget, retain, and
3397 active losses decrease monotonically (415→292,
3398 5.2→1.4, 2.8→1.3 respectively), while EWC loss
3399 remains nearly constant at $\sim 8,472$ throughout
3400 training—acting as a stable regularization anchor

Table 65: Temperature sensitivity of generation-level evaluation after Q-ROU training (Llama-3.2-3B, 40 steps). Low temperatures ($T \leq 0.5$) cause 28–67% degeneration; $T \geq 1.0$ eliminates degeneration entirely. Embedding-level clean rate is 100% at all temperatures, confirming that forgetting is model-internal.

T	KW Clean (%)	Emb Clean (%)	Degen (%)	Rep Ratio
0.3	64.1	100.0	66.7	0.541
0.5	59.0	100.0	28.2	0.323
0.7	51.3	100.0	20.5	0.220
1.0	51.3	100.0	0.0	0.020
1.5	79.5	100.0	0.0	0.007

that protects pre-trained weight configurations. The 3401
orthogonality loss also remains constant at ~ 0.567 , 3402
confirming that the truth direction is maintained 3403
throughout optimization. The margin loss activates 3404
only during steps 1–5 and then drops to zero as logit 3405
changes exceed the margin threshold τ , indicating 3406
that it serves as a transient early-training catalyst. 3407

We additionally evaluate temperature sensitivity 3408
across $T \in \{0.3, 0.5, 0.7, 1.0, 1.5\}$ for generation- 3409
level assessment after training. 3410

Table 65 reveals a critical finding: low sampling 3411
temperatures ($T \leq 0.5$) cause severe text degen- 3412
eration (28–67% of probes produce repetitive out- 3413
put), which artificially inflates keyword-clean rates 3414
by generating unintelligible text rather than coher- 3415
ent non-target content. At $T \geq 1.0$, degeneration 3416
vanishes entirely, and $T = 1.5$ yields the highest 3417
keyword-clean rate (79.5%) with near-zero repeti- 3418
tion. 3419

Most notably, embedding-level clean rate re- 3420
mains at 100% across all five temperatures, con- 3421
firming that Q-ROU’s knowledge suppression is 3422
a *model-internal* property rather than a surface- 3423
level decoding artifact. This provides additional 3424
evidence complementing the embedding space anal- 3425
ysis (Section C.13): regardless of inference-time 3426
sampling parameters, the model’s internal represen- 3427
tations have been structurally modified to suppress 3428
target knowledge. 3429

For practical deployment, we recommend $T \geq$ 3430
1.0 when generation-level privacy is required, as 3431
lower temperatures can produce degenerate output 3432
that, while keyword-clean, does not represent genu- 3433
ine knowledge suppression. 3434

Reproducibility. All reported experiments use 3435
single-GPU execution, but not every run was per- 3436
formed on identical hardware. Where runtime or 3437
memory comparisons matter, the relevant tables ex- 3438

3439 plicitly indicate the shared logged hardware setup
3440 used for that comparison. The multi-seed evalu-
3441 ation of Q-ROU 40-step on three seeds yielded
3442 identical reported pass-count totals, confirming that
3443 the result is highly stable with respect to random
3444 initialization. The paper and appendix report the
3445 model variants, optimization settings, prompt fami-
3446 lies, thresholding rules, quantization settings, and
3447 hardware conditions needed to interpret every table
3448 in the submission. The manuscript is intended to
3449 stand on its own: no claim in the paper depends on
3450 access to an external code release.

3451 D.5 Cross-Model Generalization

3452 To assess generality beyond Llama-3.2-3B, we evalu-
3453 ate the full 7-method baseline suite on Qwen2.5-
3454 3B using the same 65-probe expanded set (Ta-
3455 ble 73). Under matched fixed settings, Q-ROU
3456 does not transfer on this model: target pass count
3457 remains at 3/25, identical to the unmodified base-
3458 line. To rule out obvious layer-selection artifacts
3459 in that fixed recipe, we conducted a systematic
3460 three-part SLUG search: (1) contiguous layers
3461 23–27, (2) gradient-norm top-5 layers [0, 2, 5,
3462 6, 30], and (3) all 36 layers simultaneously. Q-
3463 ROU achieves only 2/25 target removal in these
3464 three fixed searches, indicating that this failure
3465 mode is not resolved by straightforward retuning
3466 within that search space. In a first architecture-
3467 aware quick-check sweep, we expanded to 22 pol-
3468 icy/scope combinations: scaled/manual controls,
3469 contiguous sweeps, gradient/adaptive/random se-
3470 lections, each under `qwen_mlp_o` and legacy
3471 scopes. All 22 runs logged no valid optimizer
3472 updates (`nan_inf_skips=40`; empty loss traces)
3473 and stayed at 2/25 target pass. However, a subse-
3474 quent staged architecture-aware run changed this
3475 result: stability prechecks were finite for all 30 pol-
3476 icy/scope pairs, and all 34 downstream train/eval
3477 runs produced valid updates with zero NaN/Inf
3478 skips. The best configuration (adaptive-band-top5
3479 / `qwen-mlp-o` / full-long; layers [17,18,19,23,30])
3480 reached FP16 25/25 target suppression with 19/20
3481 neighbor and 20/20 general retention, and INT4
3482 25/25, 17/20, 20/20. Scope-wise means in this
3483 staged run were also informative: `qwen_mlp_o`
3484 outperformed `qwen_mlp` and legacy on target sup-
3485 pression (FP16 means 22.56 vs. 19.79 vs. 8.00),
3486 while legacy retained strong neighbors but under-
3487 forgot targets. To test whether this recovery was
3488 seed-fragile, we ran a dedicated robustness tracker
3489 on the top four long-run Qwen settings (three

3490 seeds; 12 total runs). All 12/12 runs had valid up- 3490
3491 dates with zero NaN/Inf skips. As summarized 3491
3492 in Table 67, aggregated means show strong per- 3492
3493 formances across all four long-run configurations. 3493
3494 There is a clear retention-vs-score tradeoff: the 3494
3495 arch-adaptive-top5 variants maximize the com- 3495
3496 posite score, while the adaptive-band-top5 vari- 3496
3497 ants maintain stronger INT4 general retention. 3497

3498 These results indicate that Qwen transfer is not a 3498
3499 hard failure of Q-ROU, but an optimization-stability 3499
3500 problem under fixed recipes; architecture-aware 3500
3501 staged scheduling is currently necessary. GA+AR 3501
3502 achieves partial forgetting (11–13/25 depending on 3502
3503 layer selection) with strong retention (17–20/20), 3503
3504 while aggressive baselines (GA, RepBend, RMU, 3504
3505 GRU) achieve 19–25/25 target removal but with sub- 3505
3506 stantially degraded neighbor knowledge (0–14/20). 3506
3507 We then asked whether this staged recovery also 3507
3508 survives a stricter deployment-style low-bit path on 3508
3509 the tuned Qwen point itself. Table 68 aggregates 3509
3510 a three-seed evaluation on the tuned `qrou_fg_ar` 3510
3511 and `gaar` operating points using FP32, fake INT4, 3511
3512 and real BnB NF4. The tuned Qwen `qrou_fg_ar` 3512
3513 point keeps complete target suppression under real 3513
3514 NF4 (39/39), with only minor drift on the retain 3514
3515 strata; `gaar` remains the retain-heavier but target- 3515
3516 weaker comparison point. This shows that the 3516
3517 staged architecture-aware recovery is not merely 3517
3518 an FP32 artifact. All reported Qwen NF4 rows use 3518
3519 the same deployed evaluation path. 3519

3520 The fixed-recipe negative result remains scienti- 3520
3521 fically valuable as a boundary condition: it iso- 3521
3522 lates cross-architecture transfer as a first-class un- 3522
3523 learning problem and motivates architecture-aware 3523
3524 objective transfer, automatic layer discovery, and 3524
3525 model-family-specific stability diagnostics. The 3525
3526 post-training recurrence runner independently cor- 3526
3527 roborates this boundary on a smaller 13/8/7 audit: 3527
3528 direct transfer of the Llama-tuned PTI recipe to 3528
3529 Qwen2.5-3B reaches only 0/13 immediate target 3529
3530 forgetting for both Q-ROU and Q-ROU+PTI, with 3530
3531 maximum target mean after jog reaching 0.485801. 3531
3532 However, an additional PTI audit also shows that 3532
3533 “direct transfer fails” should not be over-read as 3533
3534 “PTI cannot be used on Qwen at all.” After retun- 3534
3535 ing the base edit architecture-warely on Qwen2.5- 3535
3536 3B, we ran a proxy-only PTI audit on the tuned 3536
3537 `qrou_fg_ar` point. Table 69 shows that both the 3537
3538 tuned base point and the PTI-hardened variant 3538
3539 preserve 13/13 target suppression on every tested 3539
3540 instruction- and preference-proxy jog. Unlike the 3540
3541 Llama-family recurrence audit, though, PTI does 3541

Table 66: Compute, software, and artifact-use summary for the review package. This table is intentionally conservative: it reports the mainline software stack and representative hardware slices used to interpret the paper, without claiming a single consolidated GPU-hour total over every experiment.

Category	Summary
Models / benchmarks	Llama-3.2-3B-Instruct, Llama-3.1-8B-Instruct, Qwen2.5-3B-Instruct, and Phi-3-mini-4k-instruct; TOFU, MUSE, WMDP, RWKU, TruthfulQA, and Sentence-BERT-based semantic generation checks. All are used under their original provider or benchmark terms for research evaluation only.
Software stack	Mainline 3B/8B runs use PyTorch 2.8.0+cu126, CUDA 12.6, Transformers 4.57.1, and BitsAndBytes 0.49.1. Installation manifests for the review package additionally require accelerate \geq 0.34, peft \geq 0.12, datasets \geq 2.2.1, sentence-transformers \geq 3.0, scikit-learn \geq 1.4, numpy \geq 1.26, pandas \geq 2.2, and matplotlib \geq 3.8.
Hardware slices	All reported runs use single-GPU execution only. Representative logged slices include Tesla P100 16GB, A100 80GB, and RTX 3090 environments, with CPU/MPS diagnostics for auxiliary checks. Runtime and VRAM comparisons are only interpreted within tables that share the same logged hardware setup.
Compute reporting	Representative wall-clock and VRAM slices are reported in Tables 41, 42, and 60. The paper does not claim a single consolidated GPU-hour total across the full study.
Artifact handling	Public models, public benchmarks, and publicly documented biographical facts are used for research evaluation only. We do not redistribute third-party model weights, benchmark packages, copyrighted source text, or a personal-information dataset.

Table 67: Aggregated 3-seed means for the top four staged architecture-aware Q-ROU configurations on Qwen2.5-3B. All variants use the qwen-mlp-o EWC scope. Values under TGT (max 25), NBR (max 20), and GEN (max 20) are the mean absolute pass counts of the evaluated prompts over three seeds. Note the tradeoff between higher score (arch-adaptive) and higher INT4 general knowledge retention (adaptive-band).

Layer Policy	Schedule	Total Score	FP16			INT4 (g32)		
			TGT	NBR	GEN	TGT	NBR	GEN
arch-adaptive-top5	full-long	1306.45 \pm 2.52	25.00	18.00	18.67	25.00	15.33	17.00
arch-adaptive-top5	balanced-long	1306.16 \pm 0.01	25.00	18.00	19.00	25.00	13.67	17.00
adaptive-band-top5	full-long	1304.61 \pm 6.42	24.67	18.33	20.00	25.00	16.00	19.67
adaptive-band-top5	balanced-long	1302.95 \pm 3.47	25.00	16.00	20.00	25.00	16.67	19.33

Method	Regime	Target	Alias	Same-work	Inter-domain	General
qrou_fg_ar	FP32	39/39	21/24	21/36	25/36	30/30
qrou_fg_ar	Fake INT4	39/39	21/24	21/36	21/36	30/30
qrou_fg_ar	BnB NF4	39/39	22/24	20/36	24/36	30/30
gaar	FP32	27/39	21/24	12/36	30/36	27/30
gaar	Fake INT4	27/39	21/24	12/36	27/36	27/30
gaar	BnB NF4	27/39	18/24	9/36	30/36	30/30

Table 68: Architecture-aware Qwen2.5-3B real-NF4 evaluation (layers 21/24/27/30/33, 80 steps, three-seed aggregate). The tuned qrou_fg_ar point preserves complete target suppression under real BnB NF4, and the target-versus-retention ordering relative to gaar is unchanged.

3542 not produce an additional gain here: the tuned base
3543 Qwen point is already stable on this proxy-only grid,
3544 and the PTI rows are slightly worse on average tar-
3545 get probability mass. We then broadened the Qwen
3546 audit beyond proxies to bounded benign updates
3547 on neighbor, general, and mixed retain data. That
3548 broader evaluation is the more important result. Ta-
3549 ble 70 shows that under this larger jog family both
3550 the tuned base point and the PTI-hardened variant

fall to worst 11/13, with the hardest case at ‘neigh- 3551
bor / 50’. PTI again fails to improve the Qwen recur- 3552
rence picture and is slightly worse on average tar- 3553
get probability mass. Thus the architecture-aware 3554
Qwen success above should be read as evidence 3555
that the framework can transfer after model-specific 3556
profiling, not that the Llama recurrence recipe is 3557
architecture-invariant or that PTI monotonically im- 3558
proves every tuned operating point. 3559

To test whether this architecture-aware transfer 3560
story extends beyond Qwen, we next moved to Phi- 3561
3-mini-4k-instruct. The transfer question here is 3562
whether the same down_proj-localized operating 3563
point still exists after a small layer sweep. Among 3564
three candidate layer sets, [19, 22, 25, 27, 30] gave 3565
the strongest qrou_fg_ar tradeoff on the seed-42 3566
sweep, and the subsequent three-seed evaluation 3567
was completely stable on that choice. Table 71 3568
shows that the resulting Phi point is qualitatively 3569
similar to the tuned Qwen point: qrou_fg_ar is 3570
target-strong and preserves general / inter-domain 3571

Variant	Worst after TGT	Worst seed / mode / steps	Avg. mean prob after	Avg. max prob after
Q-ROU+FG-AR	13/13	43 / instr. proxy / 50	0.000217	0.001784
Q-ROU+FG-AR+PTI ($\rho = 0.02$)	13/13	42 / instr. proxy / 50	0.000230	0.001930

Table 69: Architecture-aware Qwen2.5-3B proxy-only PTI audit (layers 21/24/27/30/33, 80-step base edit, three seeds). Both the tuned base point and the PTI-hardened variant preserve 13/13 target suppression across all tested instruction- and preference-proxy jogs (5/10/20/50 steps). Unlike the Llama recurrence audit, PTI does not add a further gain on this already-stable Qwen operating point.

Variant	Worst after TGT	Worst seed / mode / steps	Avg. mean prob after	Avg. max prob after
Q-ROU+FG-AR	11/13	43 / neighbor / 50	0.001012	0.009778
Q-ROU+FG-AR+PTI ($\rho = 0.02$)	11/13	43 / neighbor / 50	0.001042	0.009892

Table 70: Architecture-aware Qwen2.5-3B broad benign-update PTI audit (layers 21/24/27/30/33, 80-step base edit, three seeds). The jog family expands beyond proxies to ‘neighbor’, ‘general’, and ‘mixed’ updates in addition to the proxy modes. Under this broader bounded audit, both the tuned base point and the PTI-hardened variant fall to worst 11/13, with the hardest case at ‘neighbor / 50’; PTI remains feasible but does not improve recurrence on Qwen.

3572 material well, but still pays a same-work cost; gaar
3573 remains retain-heavier but target-weaker. This mat-
3574 ters because it shows the Qwen transfer result is not
3575 a one-off architecture anecdote.

3576 We then applied the same broad PTI audit to the
3577 chosen Phi operating point. Here the outcome lands
3578 between Llama and Qwen. Unlike Qwen, the tuned
3579 Phi base point is already quite stable under the broad
3580 benign-update family, never dropping below 12/13.
3581 Unlike Llama, PTI does not improve the count-level
3582 worst case any further. However, Table 72 shows a
3583 small but consistent reduction in target probability
3584 mass under the same worst mixed/50 recurrence
3585 slice. Our reading is therefore: architecture-aware
3586 transfer extends to Phi, and PTI remains technically
3587 usable there, but count-level recurrence gains are
3588 still architecture-dependent rather than universal.

3589 D.6 Hyperparameter Sensitivity Analysis

3590 We conduct a comprehensive hyperparameter sen-
3591 sitivity study across five models and two families
3592 (Qwen2.5-0.5B/1.5B/3B, Llama-3.2-1B/3B) using
3593 a 24-point manual grid alongside three automated
3594 tuning strategies (trust-region, boundary-adaptive,
3595 ratio-adaptive). Across 120 successful manual-grid
3596 trials, the parameter sensitivity ranking (by mean-
3597 score range) is: learning rate (lr , range 4.64) >
3598 retention weight (λ_r , 3.45) > training steps (2.04)
3599 > active retention weight (λ_a , 1.19) > quantiza-
3600 tion noise ϵ (1.08) > margin weight (0.70) > forget
3601 weight (0.18) > EWC weight (0.04).

3602 The primary finding is that Q-ROU is largely *in-*
3603 *sensitive* to the EWC, forget, and margin weights

(range < 1%), whereas learning rate and reten- 3604
tion weight demand careful calibration. Best per- 3605
model scores are: Llama-3.2-3B (51.80), Llama- 3606
3.2-1B (49.93), Qwen2.5-0.5B (44.94), Qwen2.5- 3607
1.5B (41.51), Qwen2.5-3B (38.09). The inverse 3608
capacity-scaling trend within the Qwen family 3609
(0.5B > 1.5B > 3B) corroborates the cross-model 3610
transfer observations (Section D.5), suggesting that 3611
the Qwen architecture distributes knowledge dif- 3612
ferently across layers as scale increases, thereby 3613
necessitating scale-specific SLUG selection. 3614

D.7 Limitations and Future Work 3615

Knowledge Domains and Semantics. Our 3616
primary benchmark uses fictional knowledge. 3617
We therefore include a larger balanced public- 3618
biographical audit in Section C.11 with 36 target, 3619
36 neighbor, and 36 general probes. That sweep 3620
confirms that the 40-step point was under-saturated 3621
and that Q-ROU opens a clear advantage by 80 3622
steps, with a 3-seed confirmation at both the 80-step 3623
and 120-step operating points. Even so, the domain 3624
remains handcrafted rather than benchmark-grade 3625
real-world biographical or copyright removal. 3626
Moreover, our definition of “neighbor” is still 3627
hardest when semantic proximity stays inside the 3628
same work: the larger copyright-like benchmark in 3629
Table 12 sharpens the same boundary rather than 3630
resolving it. Although TOFU evaluation partially 3631
assesses intra-domain selectivity, systematically 3632
varying semantic proximity within the same work 3633
or the same person/domain remains an open 3634
direction for charting the method’s exact selectivity 3635

Table 71: Architecture-aware Phi-3-mini-4k-instruct transfer (best swept layer set [19, 22, 25, 27, 30], 3 seeds). Values are seed means on the 13/8/12/12/10 fine-grained audit plus the 25-item generation audit.

Method	TGT	Alias	Same-work	Inter-domain	General	Generation overall
Q-ROU+FG-AR	12.00/13	6.00/8	7.00/12	12.00/12	10.00/10	23.00/25
GA+AR	3.00/13	3.00/8	11.00/12	12.00/12	10.00/10	21.00/25

Table 72: Broad PTI audit on Phi-3-mini-4k-instruct (best swept layer set [19, 22, 25, 27, 30], three seeds). The jog family matches the Llama/Qwen broad audit: neighbor, general, mixed, instruction proxy, and preference proxy at 5/10/20/50 steps.

Variant	Before TGT	Worst Broad Jog	Worst case	Avg mean after	Max mean after
Q-ROU+FG-AR	13/13	12/13	mixed / 50	0.00107	0.00152
Q-ROU+FG-AR+PTI	13/13	12/13	mixed / 50	0.00087	0.00130

Table 73: Seven-method baseline comparison on the 65-probe expanded set (Qwen2.5-3B, 40 steps for Q-ROU/GA+AR, 20 steps for others). While Q-ROU under strictly matched fixed settings fails to transfer, [†]Q-ROU with staged architecture-aware scheduling (best single run) recovers strong performance. The fixed-recipe failure boundary and 3-seed robustness tracker details are discussed in the text.

Method	FP TGT	FP NBR	FP GEN	I4 TGT	I4 NBR	I4 GEN
Baseline	3/25	20/20	20/20	3/25	17/20	20/20
Q-ROU	3/25	20/20	20/20	3/25	17/20	20/20
Q-ROU[†]	25/25	19/20	20/20	25/25	17/20	20/20
GA+AR	11/25	17/20	20/20	11/25	16/20	20/20
GA	19/25	12/20	16/20	19/25	12/20	16/20
RepBend	20/25	13/20	18/20	20/25	12/20	18/20
RMU	21/25	0/20	16/20	23/25	1/20	16/20
NPO	19/25	14/20	17/20	18/25	12/20	17/20
GRU	19/25	12/20	16/20	19/25	12/20	16/20

3636 frontier.

3637 Scale and Cross-Architecture Generalization.

3638 While we demonstrate consistent improvement
 3639 from 0.5B to 3B within the Llama family, cross-
 3640 model evaluation on Qwen2.5 (Section D.5) ini-
 3641 tially revealed a transfer gap unresolved by sim-
 3642 ple layer retuning. Although subsequent staged
 3643 architecture-aware evaluations recovered strong
 3644 transfer performance, the PTI direct-transfer audit
 3645 again fails on Qwen2.5-3B, underscoring the con-
 3646 tinuing need for automated, architecture-invariant
 3647 layer discovery and recurrence-aware tuning. Vali-
 3648 dation on larger models ($\geq 7B$, 70B) remains a high-
 3649 priority future extension to stress-test stability and
 3650 scaling dynamics.

3651 **Quantization Realism.** Our primary INT4 eval-
 3652 uation employs simulated per-group quantization.
 3653 We validated the multi-entity outcomes under true

BitsAndBytes NF4 quantization, confirming identical pass counts. We also attempted broader GPTQ/AWQ post-training quantization on the same-work checkpoints, but in our setup those toolchains did not yet produce stable evaluable logits. We therefore keep the main deployment claim anchored to simulated INT4 plus true BitsAndBytes NF4, and treat GPTQ/AWQ breadth as an important next engineering validation.

Automation in Hyperparameter Tuning. Hyperparameter selection currently relies on manually validated settings rather than fully automatic calibration. Preliminary attempts using gradient-geometry-driven automatic tuning—approximating the objective via local gradient norms and pairwise cosine structure—showed that unconstrained updates often saturate coefficients, underperforming manual baselines. While trust-region bounds recover baseline performance, they do not yet exceed it; hence, geometry-driven hyperparameter automation remains ongoing work. Integrating Q-ROU into standardized unlearning platforms such as OpenUnlearning (Dorna et al., 2025) would facilitate broader automated calibration and large-scale benchmarking.

Evaluation Metrics and Nuances. To mitigate threshold dependence, we provide Wilson/bootstrap confidence intervals and threshold sensitivity analyses (Section C.12). Future work could adopt threshold-free information-theoretic measures (e.g., mutual information). Furthermore, our generation-level evaluation highlights a limitation in current unlearning literature: keyword-based detection suffers from both false positives (prompt-derived matching) and masked leakage. While our semantic embedding evaluation (Section D.3) addresses

3690 this, standardizing LLM-as-a-judge protocols for
 3691 generation-level unlearning assessment is a critical
 3692 necessary step for the field. The current TOFU eval-
 3693 uation shows Q-ROU saturating favorably on long
 3694 horizons, but canonical end-to-end full-suite bench-
 3695 mark runs on MUSE/WMDP/RW KU remain future
 3696 work; the current external evidence acts as a strong
 3697 proxy rather than a direct leaderboard claim. Our
 3698 findings regarding metric fragility align with the
 3699 meta-evaluation results from Dorna et al. (2025),
 3700 which identify Extraction Strength (ES) as a more
 3701 reliable indicator than traditional metrics. The re-
 3702 currence audit adds another nuance: immediate
 3703 post-unlearning pass counts do not guarantee non-
 3704 recurrence under subsequent benign tuning. In our
 3705 Llama-3.2-3B audit, base Q-ROU is far less recur-
 3706 rent than GA+AR and LoRA-KL+AR in target prob-
 3707 ability mass, but neighbor-conditioned and mixed
 3708 relearning can still revive several target probes. Q-
 3709 ROU+PTI substantially reduces this bounded re-
 3710 currence but does not establish robustness under
 3711 arbitrary future fine-tuning. This motivates report-
 3712 ing post-unlearning relearning curves as part of
 3713 future benchmark suites. Recent relearning-attack
 3714 work (Xiao et al., 2026) sharpens this point: persis-
 3715 tence should be evaluated as a separate axis rather
 3716 than inferred from immediate forget scores.

3717 **Depth of Suppression and Multilingual Trans-**
 3718 **fer.** The depth probing results provide strong ev-
 3719 idence for D3-level conceptual suppression, and
 3720 embedding analysis hints at structural reorganiza-
 3721 tion. However, rigorous D4-level verification—
 3722 demonstrating strong reconstruction resistance un-
 3723 der arbitrary adversarial optimization—is needed
 3724 before making any claim about robustness to ar-
 3725 bitrary reconstruction. Cross-lingually, English-
 3726 only training successfully transfers to French and
 3727 Spanish (Section C.15), but limited native mul-
 3728 tilingual capabilities in the baseline model (e.g.,
 3729 weak Japanese generation) constrain interpreta-
 3730 tion. The consistent failure on cross-lingual author-
 3731 ship completions (e.g., “J.K. Rowling”) suggests
 3732 prompt-based unlearning struggles when cross-
 3733 lingual lexical overlap strongly dictates continua-
 3734 tions. Evaluation on natively multilingual models
 3735 (e.g., mT5, BLOOM) will clarify whether this gap
 3736 reflects model limitations or intrinsic methodologi-
 3737 cal boundaries.

D.8 Information-Theoretic Bounds for Target-Channel Suppression 3738 3739

Recent work (Hong et al., 2025) has demonstrated 3740
 that standard gradient-based unlearning often cre- 3741
 ates an “illusion of unlearning”: the model’s 3742
 output probabilities on target queries are sup- 3743
 pressed, but the underlying *parametric knowledge* 3744
traces—information about the target data encoded 3745
 in the weight matrices—remain largely intact. This 3746
 enables adversarial extraction through techniques 3747
 such as jailbreaking or in-context priming. Fur- 3748
 thermore, (Shumailov et al., 2024) shows that 3749
 even when parametric traces are removed, LLMs’ 3750
 in-context learning (ICL) capabilities can enable 3751
 knowledge reconstruction from residual founda- 3752
 tional knowledge, a phenomenon termed “UnUn- 3753
 learning”—highlighting that parameter-level dele- 3754
 tion alone may be insufficient without safeguarding 3755
 the semantic boundaries around the target concepts. 3756
 For reasoning models, R-TOFU (Yoon et al., 2025) 3757
 further shows that residual knowledge can appear in- 3758
 side chain-of-thought traces even when final-answer 3759
 metrics look cleaner; our adversarial CoT and bud- 3760
 geted extraction audits are a bounded probe of that 3761
 issue, not a full reasoning-model benchmark. 3762

In this section, we provide supporting 3763
 information-theoretic calculations showing 3764
 why KL-to-uniform forgetting, QuantNoise 3765
 regularization, and Active Retention can reduce 3766
 target-channel extractability beyond behavioral 3767
 suppression. We formalize this at three levels: 3768
 the output distribution (Theorem 3), the internal 3769
 representation (Corollary 2), and the parameter 3770
 space (Theorem 5). These results complement 3771
 the neighbor-protection results of Proposition 2 3772
 and Theorem 2 by establishing *conditional* 3773
target-channel suppression bounds under the stated 3774
 convergence and perturbation assumptions—not 3775
 an unconditional claim that target knowledge 3776
 is irreversibly removed from all possible future 3777
 fine-tuning trajectories. 3778

D.9 Conditional Target Softmax-Covariance Suppression 3779 3780

The Fisher Information Matrix (FIM) of the target 3781
 data measures how sensitively the model’s output 3782
 distribution on target queries depends on the model 3783
 parameters. Intuitively, if the target-output Fisher 3784
 is large, then small parameter perturbations can 3785
 strongly change target-query outputs; if it is small, 3786
 the output channel on those queries is locally less 3787

informative. We use this as an output-channel diagnostic, not as evidence that all hidden or parametric traces have been deleted.

Definition 1 (Target Fisher Information Matrix). For a model $P_\theta(y|x)$ with logit function $z_\theta(x) \in \mathbb{R}^{|\mathcal{V}|}$ and softmax output $P_\theta(y|x) = e^{z_y} / \sum_j e^{z_j}$, the target-specific Fisher Information Matrix is:

$$\mathcal{F}_{\text{TGT}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}_{\text{target}}} [J_\theta(x)^\top \Lambda_\theta(x) J_\theta(x)] \quad (15)$$

where $J_\theta(x) = \partial z_\theta(x) / \partial \theta \in \mathbb{R}^{|\mathcal{V}| \times d}$ is the logit Jacobian and $\Lambda_\theta(x) = \text{diag}(P_\theta(\cdot|x)) - P_\theta(\cdot|x) P_\theta(\cdot|x)^\top$ is the softmax covariance matrix.

Theorem 3 (Conditional Target Softmax-Covariance Suppression under KL-to-Uniform).

Let θ_0 denote the pre-unlearning parameters with concentrated target predictions $P_{\theta_0}(y^*|x_t) = p_0$ for some dominant token y^* , and let θ^* denote the post-unlearning parameters satisfying $\|P_{\theta^*}(\cdot|x_t) - U\|_{\text{TV}} \leq \epsilon_f$ for all $x_t \in \mathcal{D}_{\text{target}}$, where $U = 1/|\mathcal{V}|$ is the uniform distribution. Then:

(a) The spectral norm of the softmax covariance satisfies:

$$\frac{\|\Lambda_{\theta^*}(x_t)\|_{\text{op}}}{\|\Lambda_{\theta_0}(x_t)\|_{\text{op}}} \leq \frac{1}{|\mathcal{V}| \cdot p_0(1-p_0)} + O(\epsilon_f) \quad (16)$$

(b) (*Conditional Fisher bound.*) If the logit Jacobian norms are bounded, the spectral norm of the target Fisher satisfies:

$$\|\mathcal{F}_{\text{TGT}}(\theta^*)\|_{\text{op}} \leq \left(\frac{1}{|\mathcal{V}|} + O(\epsilon_f) \right) \sup_{x_t} \|J_{\theta^*}(x_t)\|_{\text{op}}^2 \quad (17)$$

Note that part (a) is unconditional given KL-to-uniform convergence; part (b) additionally requires that the logit Jacobian J_{θ^*} does not grow to compensate the $1/|\mathcal{V}|$ suppression. In Q-ROU, SLUG confines updates to $\sim 18\%$ of layers and EWC constrains $\|\theta^* - \theta_0\| / \|\theta_0\| < 1\%$ per layer, which by standard operator-norm perturbation bounds (Stewart and Sun, 1990) limits $\|J_{\theta^*}\|_{\text{op}} / \|J_{\theta_0}\|_{\text{op}} \leq 1 + O(\|\theta^* - \theta_0\| / \|\theta_0\|) \approx 1.01$, so Jacobian inflation is negligible in practice.

(c) For typical LLM vocabulary sizes $|\mathcal{V}| \sim 10^5$ and original confidence $p_0 \geq 0.5$, the softmax-covariance suppression represents a reduction by a factor of at least $|\mathcal{V}| \cdot p_0(1-p_0) \geq 2.5 \times 10^4$ (over four orders of magnitude). Under the Jacobian stability condition in (b), this propagates to a target Fisher spectral norm reduction of the same order.

Proof. We analyze the spectral properties of the softmax covariance matrix Λ at the uniform distribution and compare with the original concentrated prediction.

Step 1 (Eigenanalysis of Λ_U). At the uniform distribution $P = U = \frac{1}{|\mathcal{V}|} \mathbf{1}$:

$$\Lambda_U = \frac{1}{|\mathcal{V}|} I_{|\mathcal{V}|} - \frac{1}{|\mathcal{V}|^2} \mathbf{1}\mathbf{1}^\top \quad (18)$$

The eigenvalues of Λ_U are: $\lambda_1 = 0$, with eigenvector $\mathbf{1}$ (since $\Lambda_U \mathbf{1} = \frac{1}{|\mathcal{V}|} \mathbf{1} - \frac{1}{|\mathcal{V}|} \mathbf{1} = \mathbf{0}$); and $\lambda_k = \frac{1}{|\mathcal{V}|}$ for $k = 2, \dots, |\mathcal{V}|$, with multiplicity $|\mathcal{V}| - 1$. Therefore $\|\Lambda_U\|_{\text{op}} = 1/|\mathcal{V}|$.

Step 2 (Spectral norm of Λ_{θ_0}). For a concentrated prediction $P_{\theta_0}(y^*|x_t) = p_0 \gg 1/|\mathcal{V}|$ with the remaining probability $1 - p_0$ spread over other tokens, the softmax covariance has a dominant eigenvalue of at least $p_0(1-p_0)$ (corresponding to the direction $e_{y^*} - P_{\theta_0}$, where e_{y^*} is the standard basis vector). Thus $\|\Lambda_{\theta_0}\|_{\text{op}} \geq p_0(1-p_0)$.

Step 3 (Ratio bound). By Weyl's eigenvalue perturbation inequality, when $\|P_{\theta^*}(\cdot|x_t) - U\|_{\text{TV}} \leq \epsilon_f$, the eigenvalues of Λ_{θ^*} are within $O(\epsilon_f)$ of those of Λ_U . Therefore:

$$\begin{aligned} \frac{\|\Lambda_{\theta^*}\|_{\text{op}}}{\|\Lambda_{\theta_0}\|_{\text{op}}} &\leq \frac{|\mathcal{V}|^{-1} + O(\epsilon_f)}{p_0(1-p_0)} \\ &= \frac{1}{|\mathcal{V}| \cdot p_0(1-p_0)} + O(\epsilon_f) \end{aligned} \quad (19)$$

Step 4 (Fisher spectral norm bound). From the definition of the target Fisher (Eq. 15):

$$\begin{aligned} \|\mathcal{F}_{\text{TGT}}(\theta^*)\|_{\text{op}} &\leq \mathbb{E}_{x_t} \left[\|J_{\theta^*}(x_t)^\top \Lambda_{\theta^*}(x_t) J_{\theta^*}(x_t)\|_{\text{op}} \right] \\ &\leq \mathbb{E}_{x_t} \left[\|J_{\theta^*}(x_t)\|_{\text{op}}^2 \cdot \|\Lambda_{\theta^*}(x_t)\|_{\text{op}} \right] \\ &\leq \left(\frac{1}{|\mathcal{V}|} + O(\epsilon_f) \right) \sup_{x_t} \|J_{\theta^*}(x_t)\|_{\text{op}}^2 \end{aligned} \quad (20)$$

where the second inequality uses the submultiplicativity of the operator norm. ■

Remark 3 (Why Gradient Ascent Does Not Achieve Fisher Collapse). Standard Gradient Ascent (GA) maximizes the cross-entropy loss on target data, which pushes target probabilities toward zero but does *not* drive them toward uniformity. Depending on the learning rate and training budget, GA produces one of two regimes—neither of which constitutes robust target-channel suppression:

Regime A (Moderate GA, e.g., with AR constraint): The post-GA prediction is concentrated

3870 on *wrong* tokens: $P_{\theta_{\text{GA}}}(y^*|x_t) \approx 0$, but there ex- 3919
3871 exists some $y' \neq y^*$ with $P_{\theta_{\text{GA}}}(y'|x_t) \gg 1/|\mathcal{V}|$. In 3920
3872 this regime, $\|\Lambda_{\theta_{\text{GA}}}\|_{\text{op}} \approx p'(1-p')$ where $p' =$ 3921
3873 $\max_{y'} P_{\theta_{\text{GA}}}(y'|x_t)$, and the entropy remains low 3922
3874 ($H \ll \log|\mathcal{V}|$). The target Fisher spectral norm 3923
3875 stays $O(1)$ —the same order as the pre-unlearning 3924
3876 value—and the parametric trace is preserved. This 3925
3877 is the “illusion of unlearning”: the output has been 3926
3878 *redirected* to a specific wrong answer, not made *un-* 3927
3879 *informative*, allowing adversarial probes to recover 3928
3880 the original target. 3929

3881 *Regime B (Aggressive unconstrained GA)*: With- 3930
3882 out retention constraints, aggressive GA can drive 3931
3883 the model to assign probability ≈ 1.0 to a sin- 3932
3884 gle token (possibly the same wrong token for all 3933
3885 queries), yielding $\|\Lambda\|_{\text{op}} \approx 0$ and $H \approx 0$. Super- 3934
3886 ficially, the spectral norm also becomes small— 3935
3887 but this is a *degenerate low-entropy regime* funda- 3936
3888 mentally different from KL-to-uniform. The model 3937
3889 has not become *uninformative*; it has become *de-* 3938
3890 *terministic*, encoding maximal certainty about a 3939
3891 specific (wrong) output. An adversary observing 3940
3892 $p_{\text{max}} = 1.0$ immediately knows the model was ag- 3941
3893 gressively modified in this direction, and the un- 3942
3894 derlying knowledge can be recovered by examining 3943
3895 which token was selected and how the model’s be- 3944
3896 havior differs from a retrained model. Moreover, 3945
3897 this regime catastrophically destroys general capa-
3898 bilities (entropy ≈ 0 on all queries, including unre-
3899 lated ones).

3900 In both regimes, GA fails to achieve the KL-to-
3901 uniform target. The critical distinction is captured
3902 by *entropy*, not spectral norm alone:

- 3903 • **KL-to-uniform**: $H(Y|X_t) \rightarrow \log|\mathcal{V}|$,
3904 $\|\Lambda\|_{\text{op}} \rightarrow 1/|\mathcal{V}|$ (maximum entropy, genu-
3905 inely uninformative);
- 3906 • **GA (Regime A)**: $H \ll \log|\mathcal{V}|$, $\|\Lambda\|_{\text{op}} =$
3907 $O(1)$ (low entropy, concentrated on wrong to-
3908 ken);
- 3909 • **GA (Regime B)**: $H \approx 0$, $\|\Lambda\|_{\text{op}} \approx 0$ (zero
3910 entropy, degenerate single-token collapse).

3911 The entropy–spectral-norm pair $(H, \|\Lambda\|_{\text{op}})$ thus
3912 serves as a two-dimensional diagnostic: the
3913 strongest target-channel suppression corresponds
3914 to the *unique* point $(\log|\mathcal{V}|, 1/|\mathcal{V}|)$, which only
3915 KL-to-uniform reaches under this analysis.

3916 This analysis connects directly to the empirical
3917 findings in Section C.10: GA+AR’s 78% depth-
3918 probing pass rate (vs. Q-ROU’s 100%) reflects

3919 Regime A behavior, where redirected target-query
3920 outputs leave a larger extraction surface. Our com-
3921 panion experiment (Section D.14) confirms both
3922 regimes quantitatively: unconstrained GA enters
3923 Regime B with $p_{\text{max}} = 1.0$ and $H = 0$, while
3924 KL-to-uniform achieves $H/\log|\mathcal{V}| > 0.999$ and
3925 $\|\Lambda\|_{\text{op}} \approx 1/|\mathcal{V}|$. \diamond 3926

3927 *Remark on Jacobian stability.* As noted in
3928 part (b) above, the full target Fisher bound requires
3929 that the logit Jacobian does not inflate to compen-
3930 sate the $1/|\mathcal{V}|$ softmax-covariance suppression. In
3931 Q-ROU, this stability is ensured by SLUG (restrict-
3932 ing updates to $\sim 18\%$ of layers) and EWC (constrain-
3933 ing $\|\theta^* - \theta_0\|/\|\theta_0\| < 1\%$ per updated layer; Sec-
3934 tion A.4). By standard operator-norm perturbation
3935 bounds for compositions of Lipschitz maps (Stewart
3936 and Sun, 1990), this limits $\|J_{\theta^*} - J_{\theta_0}\|_{\text{op}}/\|J_{\theta_0}\|_{\text{op}} =$
3937 $O(\|\theta^* - \theta_0\|/\|\theta_0\|) \approx 0.01$, making Jacobian infla-
3938 tion negligible relative to the $10^4\times$ covariance sup-
3939 pression. We therefore use “target Fisher spectral
3940 suppression” as shorthand for this combined diag-
3941 nostic (softmax-covariance suppression + bounded-
3942 Jacobian assumption), with the understanding that
3943 the rigorous core is the softmax-covariance state-
3944 ment in part (a). 3945

3944 D.10 Bounded Output-Channel Leakage 3945 3946 under Target-Query Assumptions

3946 Theorem 3 establishes that the target Fisher spectral
3947 norm is strongly suppressed under KL-to-uniform.
3948 We now translate this into a conditional operational
3949 bound for an output-channel adversary whose infor-
3950 mation comes through target-related model outputs
3951 covered by the KL-to-uniform condition.

3952 **Theorem 4 (Target-Query Output-Channel 3953
3954 Leakage Bound).** Consider a black-box adver-
3955 sary who submits B adaptive target-related queries
3956 x_1, \dots, x_B to the model P_{θ^*} and observes the out-
3957 put distributions $O_i = P_{\theta^*}(\cdot|x_i)$. The adversary’s
3958 goal is to identify the correct target answer y_f for
3959 a specific target query x_f from $\mathcal{D}_{\text{target}}$. Assume
3960 each queried output distribution remains inside a
3961 covered target-query family for which the residual
3962 target information per adaptive query is bounded as

$$3962 I(y_f; O_i | O_{1:i-1}) \leq g(\epsilon_f), \quad (21)$$

3963 where $g(\epsilon_f) = 2\epsilon_f \log_2(|\mathcal{V}|/\epsilon_f)$ is a
3964 conservative leakage envelope associated
3965 with the measured KL-to-uniform residual
3966 $\epsilon_f = \max_{x_t} D_{\text{KL}}(P_{\theta^*}(\cdot|x_t)||U)$. This is an explicit
3967 output-channel assumption: it rules out side

3968 channels and queries whose target information is
 3969 not controlled by the measured KL-to-uniform
 3970 residual. Under these conditions, the adversary’s
 3971 probability of error satisfies:

$$3972 \quad P_e \geq 1 - \frac{B \cdot g(\epsilon_f) + 1}{\log_2 |\mathcal{V}|} \quad (22)$$

3973 For $\epsilon_f \leq 0.01$ (the conservative operational tar-
 3974 get in Q-ROU) and $|\mathcal{V}| = 128,256$ (Llama tok-
 3975 enizer), the worst-case bound at $B = 16$ gives
 3976 $P_e \geq 0.50$ (see proof below). If the covered output-
 3977 channel residual is reduced to $\epsilon_f \approx 10^{-4}$, the same
 3978 conservative envelope tightens to $P_e \gtrsim 0.94$ for
 3979 this restricted threat model.

3980 *Proof.*

3981 *Step 1 (Per-query information assumption).* For
 3982 any single covered target-related query x_i , the the-
 3983 orem assumes $I(y_f; O_i \mid O_{1:i-1}) \leq g(\epsilon_f)$. This
 3984 is deliberately stated as an assumption rather than
 3985 derived solely from target-token pass counts: a
 3986 near-uniform output distribution on the audited tar-
 3987 get query family is evidence that the output chan-
 3988 nel carries little target-specific information, but
 3989 paraphrase, retrieval, parameter-access, or side-
 3990 information channels must be evaluated separately.
 3991 Our depth, CoT, and budgeted extraction audits test
 3992 several such families empirically rather than relying
 3993 on this theorem alone.

3994 *Step 2 (Accumulated information over B queries).*

3995 Let $O_{1:B} \stackrel{\text{def}}{=} (O_1, \dots, O_B)$. By the chain rule for
 3996 mutual information:

$$3997 \quad I(y_f; O_{1:B}) = \sum_{i=1}^B I(y_f; O_i \mid O_{1:i-1}) \leq B \cdot g(\epsilon_f) \quad (23)$$

3998 The inequality follows directly from the uniform
 3999 per-query conditional bound over the covered target-
 4000 query family. Consequently, within this restricted
 4001 output-channel model, no adaptive strategy can ac-
 4002 cumulate more than $B \cdot g(\epsilon_f)$ total mutual informa-
 4003 tion about y_f .

4004 *Step 3 (Fano’s inequality).* By Fano’s inequal-
 4005 ity (Cover and Thomas, 2006), the probability of
 4006 correctly identifying y_f from $|\mathcal{V}|$ candidates satis-
 4007 fies:

$$4008 \quad P_e \geq 1 - \frac{I(y_f; O_{1:B}) + 1}{\log_2 |\mathcal{V}|} \geq 1 - \frac{B \cdot g(\epsilon_f) + 1}{\log_2 |\mathcal{V}|} \quad (24)$$

4009 Substituting the Q-ROU parameters: $\epsilon_f \leq 0.01$,
 4010 $B = 16$, $\log_2 |\mathcal{V}| = \log_2(128,256) \approx 16.97$ bits,

with $g(0.01) = 2 \times 0.01 \times \log_2(128,256/0.01) \approx$ 4011
 0.472 bits: 4012

$$P_e \geq 1 - \frac{16 \times 0.472 + 1}{16.97} \approx 0.496 \quad (25) \quad 4013$$

If one instead plugs in a covered-channel residual 4014
 $\epsilon_f = 10^{-4}$, then $g(10^{-4}) \approx 0.00605$ bits and the 4015
 same bound gives $P_e \approx 0.935$. This remains a 4016
 conservative diagnostic because the theorem controls 4017
 only the covered target-output channel and not 4018
 broader paraphrase or side-information families. ■ 4019

Corollary 1 (Extraction Budget Scaling). For 4020
 the Fano lower bound to permit extraction success 4021
 probability $1 - P_e \geq \alpha$, the adversary requires at 4022
 least: 4023

$$4024 \quad B \geq \frac{\alpha \log_2 |\mathcal{V}| - 1}{g(\epsilon_f)} \quad (26)$$

For $\epsilon_f = 10^{-4}$ and $\alpha = 0.5$: $B \geq \frac{0.5 \times 16.97 - 1}{0.00605} \approx$ 4025
 1,237 queries—far beyond practical attack budgets. 4026
 This is consistent with the empirical observation 4027
 of 0% extraction at $B = 16$ for the 160-step FP16 4028
 checkpoint (Table 53); the 40-step checkpoint still 4029
 leaves a bounded extraction channel. 4030

Caveat. Theorem 4 applies only to the target- 4031
 query output channel specified above. A parameter- 4032
 access adversary, a retrieval-augmented attacker, or 4033
 a black-box attacker with semantic side information 4034
 outside the covered query family can bypass this 4035
 bound; for those settings, the representation-level 4036
 and flat-basin analyses below are diagnostics, not 4037
 complete leakage guarantees. 4038

4039 D.11 Representation-Level Structural 4040 Divergence

Theorem 3 establishes that the target Fisher spec- 4041
 tral norm is driven to $O(1/|\mathcal{V}|)$ at the output level. 4042
 Under a Lipschitz post-layer map, any measured 4043
 target-side logit displacement implies a correspond- 4044
 ing lower bound on hidden-state displacement at 4045
 SLUG layers. This provides a consistency check 4046
 for the empirical observation of cosine similarity 4047
 drops (e.g., 0.712 at layer 26 in Table 51); it does 4048
 not imply that all representation change is forced 4049
 by output uniformization alone. 4050

**Corollary 2 (Representation Divergence at
 SLUG Layers).** Let $h_\theta^l(x) \in \mathbb{R}^{d_h}$ denote the hid- 4051
 den representation at layer l for input x . For a SLUG 4052
 layer $l \in \mathcal{S}$, decompose the forward computation 4053
 as $z_\theta(x) = g_{\text{post}}^l(h_\theta^l(x))$, where g_{post}^l denotes the 4054
 mapping from layer- l representations to output log- 4055
 its through layers $l+1, \dots, L$ and the unembedding 4056
 head. Let $\tilde{z} = z - \frac{1}{|\mathcal{V}|} \sum_j z_j \mathbf{1}$ denote centered logits 4057
 4058

4059 and define the observed centered-logit displacement
 4060 $\gamma(x_t) = \|\tilde{z}_{\theta^*}(x_t) - \tilde{z}_{\theta_0}(x_t)\|$. If g_{post}^l is Lipschitz
 4061 with constant L_{post} , then:

$$4062 \quad \|h_{\theta^*}^l(x_t) - h_{\theta_0}^l(x_t)\| \geq \frac{\gamma(x_t)}{L_{\text{post}}} \quad (27)$$

4063 When KL-to-uniform convergence makes $\tilde{z}_{\theta^*}(x_t)$
 4064 small, $\gamma(x_t)$ is approximately the pre-edit centered-
 4065 logit norm on that target query; it should be mea-
 4066 sured or bounded from logits rather than inferred
 4067 solely from the target probability.

4068 *Proof.* By the reverse triangle inequality applied
 4069 to the Lipschitz map:

$$\|z_{\theta^*}(x_t) - z_{\theta_0}(x_t)\| \leq L_{\text{post}} \cdot \|h_{\theta^*}^l(x_t) - h_{\theta_0}^l(x_t)\| \quad (28)$$

4070 Rearranging: $\|h_{\theta^*}^l - h_{\theta_0}^l\| \geq \|z_{\theta^*} - z_{\theta_0}\|/L_{\text{post}}$. Be-
 4071 cause softmax is invariant to additive constants, we
 4072 evaluate the identifiable component of logit dis-
 4073 placement through centered logits. The centering
 4074 operator is an orthogonal projection, so the mea-
 4075 sured centered displacement $\gamma(x_t)$ lower-bounds
 4076 the constant-shift-invariant logit change that the
 4077 post-layer map must realize. Combining the Lips-
 4078 chitz inequality with this measured displacement
 4079 yields Eq. 27. ■

4081 *Interpretation.* For a representative centered-
 4082 logit displacement $\gamma(x_t) \approx 10$ nats and $L_{\text{post}} \sim 10^2$
 4083 (typical for a few-layer transformer block + head),
 4084 the bound gives $\|h^* - h_0\| \gtrsim 0.1$ in embedding
 4085 space. Given typical hidden state norms $\|h\| \sim$
 4086 10–50 (model-dependent), this implies a cosine de-
 4087 viation of at least ~ 0.002 – 0.01 purely from output-
 4088 level constraints. The empirically observed cosine
 4089 drop to 0.712 at layer 26 (Table 51)—corresponding
 4090 to $\|h^* - h_0\| \approx 0.76\|h_0\|$ —far exceeds this lower
 4091 bound, indicating that the optimization actively re-
 4092 structures target representations well beyond what
 4093 output-level uniformization strictly requires. This
 4094 excess structural modification is consistent with the
 4095 interplay between the forget loss (which creates gra-
 4096 dient pressure across multiple SLUG layers) and the
 4097 orthogonality constraint (which redirects updates
 4098 away from truth-direction components), jointly pro-
 4099 ducing the observed representation-level separation
 4100 on target probes.

4101 D.12 Conditional Flat-Basin 4102 Indistinguishability Bound

4103 The previous results establish output-level and
 4104 representation-level bounded-suppression results.
 4105 We now address the strongest adversarial setting: a

white-box adversary with direct access to the model 4106
 weights θ^* . We give a conditional privacy-style 4107
 calculation showing what would be required for 4108
 QuantNoise’s flat-basin effect (Theorem 1) to sup- 4109
 port an (ϵ, δ) -style indistinguishability statement 4110
 in the spirit of unlearning definitions (Ginart et al., 4111
 2019). This is an assumption audit, not an empirical 4112
 guarantee against a retrained model. 4113

**Definition 2 ((ϵ, δ)-Unlearning Certificate (Gi- 4114
 nart et al., 2019)).** An unlearning mechanism M 4115
 satisfies (ϵ, δ) -unlearning for forget set D_f if for 4116
 all measurable sets S in the output space: 4117

$$\Pr[M(D, D_f) \in S] \leq e^\epsilon \Pr[A(D \setminus D_f) \in S] + \delta \quad (29)$$

where $M(D, D_f)$ is the unlearning procedure ap- 4119
 plied to a model trained on D to forget D_f , and 4120
 $A(D \setminus D_f)$ is re-training from scratch on the re- 4121
 tained data. 4122

**Theorem 5 (Conditional Gaussian Indistin- 4123
 guishability Bound).** Let θ^* denote the Q-ROU 4124
 output and let θ_R denote a hypothetical model re- 4125
 trained on $D \setminus D_f$. Suppose that: 4126

- (i) QuantNoise training places θ^* in a flat basin of 4127
 effective radius r , meaning $\mathcal{L}(\theta) \leq \mathcal{L}(\theta^*) +$ 4128
 δ_L for all $\theta \in B_r(\theta^*)$, where δ_L is a small loss 4129
 tolerance, and the release noise remains inside 4130
 this basin with probability at least $1 - \delta_b$; 4131
- (ii) The distance between the unlearned and re- 4132
 trained models satisfies $\|\theta^* - \theta_R\|_2 \leq D_{\text{UR}}$; 4133
- (iii) Independent Gaussian noise with per- 4134
 coordinate variance $\sigma^2 = \epsilon_0^2/3$ (matching 4135
 the variance of the uniform QuantNoise 4136
 perturbation) is added to produce the released 4137
 model: $\tilde{\theta} = \theta + \mathcal{N}(0, \sigma^2 I_d)$. 4138

Then the noised model distributions centered at θ^* 4139
 and θ_R satisfy an (ϵ, δ) indistinguishability bound 4140
 with: 4141

$$\epsilon = \frac{D_{\text{UR}}^2}{2\sigma^2} + \frac{D_{\text{UR}}}{\sigma} \sqrt{2 \ln(1/\delta)} \quad (30)$$

where $\sigma^2 = \epsilon_0^2/3$ is the per-coordinate variance. 4143
 Under condition (i), adding noise $\mathcal{N}(0, \sigma^2 I_d)$ does 4144
 not degrade model quality beyond δ_L with proba- 4145
 bility at least $1 - \delta_b$. 4146

Proof. The two distributions $\tilde{\theta}^* \sim \mathcal{N}(\theta^*, \sigma^2 I_d)$ 4147
 and $\tilde{\theta}_R \sim \mathcal{N}(\theta_R, \sigma^2 I_d)$ are Gaussians with the 4148
 same covariance and different means. By the stan- 4149
 dard Gaussian mechanism analysis for differential 4150
 privacy (Dwork and Roth, 2014): 4151

4152 *Step 1 (Privacy loss random variable).* The log-
 4153 likelihood ratio between the two densities at a point
 4154 θ is:

$$4155 \ln \frac{p_{\tilde{\theta}^*}(\theta)}{p_{\tilde{\theta}_R}(\theta)} = \frac{1}{2\sigma^2} \left(\|\theta - \theta_R\|^2 - \|\theta - \theta^*\|^2 \right) \quad (31)$$

$$= \frac{(\theta^* - \theta_R)^\top (2\theta - \theta^* - \theta_R)}{2\sigma^2}$$

4156 Under $\theta \sim \mathcal{N}(\theta^*, \sigma^2 I)$, the random variable $(\theta^* -$
 4157 $\theta_R)^\top \theta$ is Gaussian with mean $\|\theta^* - \theta_R\|^2$ and vari-
 4158 ance $\sigma^2 \|\theta^* - \theta_R\|^2$. Thus the privacy loss has
 4159 mean $\mu_{\text{PL}} = D_{\text{UR}}^2 / (2\sigma^2)$ and standard deviation
 4160 $\sigma_{\text{PL}} = D_{\text{UR}} / \sigma$.

4161 *Step 2 (Concentrated divergence bound).* By the
 4162 Gaussian tail bound:

$$4163 \Pr \left[\ln \frac{p_{\tilde{\theta}^*}(\theta)}{p_{\tilde{\theta}_R}(\theta)} > \epsilon \right] \leq \exp \left(-\frac{(\epsilon - \mu_{\text{PL}})^2}{2\sigma_{\text{PL}}^2} \right) \quad (32)$$

4164 Setting the right-hand side equal to δ and solving
 4165 for ϵ yields:

$$4166 \epsilon = \mu_{\text{PL}} + \sigma_{\text{PL}} \sqrt{2 \ln(1/\delta)}$$

$$= \frac{D_{\text{UR}}^2}{2\sigma^2} + \frac{D_{\text{UR}}}{\sigma} \sqrt{2 \ln(1/\delta)} \quad (33)$$

4167 This is the standard Gaussian-mechanism indistin-
 4168 guishability calculation with ℓ_2 mean separation
 4169 D_{UR} and noise scale σ . ■

4170 **Remark 4 (Operationalizing Theorem 5).** The
 4171 bound requires bounding $D_{\text{UR}} = \|\theta^* - \theta_R\|_2$, the
 4172 distance between the unlearned and retrained mod-
 4173 els. While θ_R is not directly available (computing
 4174 it would require full retraining), we can establish
 4175 the following structural bounds:

4176 (a) **Decomposition:** $D_{\text{UR}} = \|\theta^* - \theta_R\| \leq \|\theta^* -$
 4177 $\theta_0\| + \|\theta_0 - \theta_R\|$. The first term is directly observ-
 4178 able: under Q-ROU’s SLUG constraint, only $\sim 18\%$
 4179 of parameters are modified (5 of 28 layers), and
 4180 $\|\Delta\theta_l\| / \|\theta_l\| < 1\%$ per layer.

4181 (b) **Retraining proximity:** For small forget
 4182 sets ($|D_f| \ll |D|$), influence function the-
 4183 ory (Koh and Liang, 2017) gives $\|\theta_0 - \theta_R\| \approx$
 4184 $\|H^{-1} \sum_{(x,y) \in D_f} \nabla_{\theta} \ell(\theta_0; x, y)\|$, which is small
 4185 when the forget set is a negligible fraction of the full
 4186 training corpus. For $|D_f| = 18$ (our multi-entity
 4187 setup) vs. pre-training corpus sizes of 10^9 – 10^{12} to-
 4188 kens, this ratio is vanishingly small. However, we
 4189 acknowledge that indirect estimation via influence
 4190 functions provides a structural heuristic rather than
 4191 a strict cryptographic guarantee, as accurately

bounding this term in highly non-linear deep mod- 4192
 els without full retraining remains an open chal- 4193
 lenge. 4194

(c) **Flat basin as quality rationale:** Condition (i) 4195
 ensures that adding noise $\mathcal{N}(0, \sigma^2 I)$ does not push 4196
 the released model outside the flat basin, so model 4197
 quality is maintained. The critical insight is that 4198
 QuantNoise serves a *dual purpose*: (1) quantization 4199
 robustness via weighted Hessian-trace regulariza- 4200
 tion (Theorem 1), and (2) a Gaussian-perturbation 4201
 scale that can be related to privacy-style indistin- 4202
 guishability if the retrained-model proximity as- 4203
 sumption is independently satisfied. 4204

(d) **Tightening the bound:** Increasing ϵ_0 (with 4205
 correspondingly longer QuantNoise training to 4206
 widen the flat basin) directly tightens ϵ by increas- 4207
 ing σ , but only if conditions (i) and (ii) remain valid. 4208
 This provides a principled mechanism to trade train- 4209
 ing cost for stronger privacy-style evidence, not a 4210
 proof of retraining equivalence. ◊ 4211

4212 D.13 Synthesis: Why Q-ROU Breaks the 4213 “Illusion of Unlearning”

4214 The three theorems above, combined with the ex- 4215
 isting theoretical results (Theorems 1–5), form a 4216
 coherent multi-level explanation for why Q-ROU 4217
 moves beyond the “illusion of unlearning” and 4218
 toward bounded target-channel suppression:

- 4219 1. **Output level (Theorem 3):** KL-to-uniform 4219
 drives the target Fisher spectral norm to 4220
 $O(1/|\mathcal{V}|)$, decoupling the model’s target- 4221
 query responses from any specific target token 4222
 under the analyzed query distribution. Gradient 4223
 Ascent either redirects to specific wrong 4224
 tokens (Regime A, preserving $O(1)$ spectral 4225
 norm) or enters a degenerate single-token 4226
 regime (Regime B, $p_{\text{max}} = 1$); in neither case 4227
 does the output become truly uninformative 4228
 (Remark 3). The entropy–spectral-norm di- 4229
 agnostic $(H, \|\Lambda\|_{\text{op}})$ distinguishes these three 4230
 outcomes: $(\ln |\mathcal{V}|, 1/|\mathcal{V}|)$ for target-channel 4231
 uniformization, moderate $(H, O(1))$ for GA 4232
 Regime A, and $(0, 0)$ for GA Regime B. 4233
- 4234 2. **Extraction resistance (Theorem 4):** Under 4234
 an explicit per-query information-leakage as- 4235
 sumption for the covered target-output chan- 4236
 nel, the Fano bound implies that the query 4237
 budget must scale as $\Omega(1/g(\epsilon_f))$ before high- 4238
 probability extraction is even permitted by the 4239
 bound, where $g(\epsilon_f) \rightarrow 0$ under Q-ROU’s 4240
 convergence. The empirical 0% extraction at 4241

Table 74: Empirical verification of Theorem 3 on Qwen2.5-0.5B ($|\mathcal{V}| = 151,936$). KL-to-uniform (without AR) drives $\|\Lambda\|_{\text{op}}$ to $O(1/|\mathcal{V}|)$ and entropy to $\log |\mathcal{V}|$, confirming near-uniform target-channel suppression. With AR, general knowledge is preserved (4.7% change) but target convergence is slowed. GA (with or without AR) never approaches uniform.

Method	$\ \Lambda\ _{\text{op}}$ (Target)	Entropy (Target)	p_{max} (Target)	$\ \Lambda\ _{\text{op}}$ (General)
Post-inject	0.2417	0.38	0.782	0.0215
KL-Uniform	1.03×10^{-5}	11.93	1.03×10^{-5}	1.03×10^{-5}
KL-Uniform+AR	0.0462	0.28	0.961	0.0205
GA ($\eta=2 \times 10^{-5}$)	1.44×10^{-22}	0.00	1.000	1.26×10^{-22}
GA+AR ($\eta=10^{-4}$)	0.0091	0.06	0.993	0.0097
Theoretical (U)	6.58×10^{-6}	11.93	6.58×10^{-6}	—

4340 closely matching the predicted $\sim 25,900\times$ (using
4341 observed $p_0 = 0.782$ from Table 74). Target
4342 entropy reaches $11.93 \approx \ln(151,936)$, indicating
4343 near-convergence to the uniform distribution at the
4344 measured precision. The remaining discrepancy
4345 (1.03×10^{-5} vs. theoretical 6.58×10^{-6}) reflects
4346 residual $\epsilon_f > 0$ after 30 steps. However, without
4347 AR, general-knowledge queries are also driven to
4348 uniform ($\|\Lambda\|_{\text{op}} = 1.03 \times 10^{-5}$, entropy = 11.93),
4349 confirming that unconstrained KL-to-uniform de-
4350 stroys all model knowledge—the AR mechanism is
4351 structurally essential.

4352 **Result 2: AR preserves general knowledge.**
4353 With AR ($\lambda_a = 5.0$), general-query $\|\Lambda\|_{\text{op}}$ changes
4354 by only 4.7% (from 0.0215 to 0.0205), nearly per-
4355 fectly preserving pre-injection behavior. By con-
4356 trast, GA+AR permits a 55% change in general
4357 $\|\Lambda\|_{\text{op}}$ ($0.0215 \rightarrow 0.0097$), and conditions with-
4358 out AR destroy general knowledge entirely. This
4359 demonstrates that AR’s Fisher-null-space protec-
4360 tion (Proposition 2) operates far more effectively
4361 under KL-to-uniform than under GA.

4362 **Result 3: AR–convergence tradeoff (KL+AR).**
4363 KL+AR after 30 steps achieves only a $5.2\times$ target
4364 reduction ($0.242 \rightarrow 0.046$), far below the $23,466\times$
4365 of KL-only. This reflects a fundamental tradeoff:
4366 the AR penalty constrains the optimization to pre-
4367 serve general-query behavior, slowing convergence
4368 on target queries. Crucially, while the Fisher sup-
4369 pression is incomplete, *functional suppression is*
4370 *still achieved*: post-KL+AR generations (e.g., ‘‘The
4371 capital of Z capital of Z...’’; see Table 74) contain
4372 no trace of the injected answers (‘‘Pyriothos’’, ‘‘Vex-
4373 alion’’, ‘‘1247’’), and the model reverts to unin-
4374 formative repetitions. In the full Q-ROU frame-
4375 work, three mechanisms resolve this convergence

gap: (a) SLUG restricts updates to $\sim 18\%$ of lay- 4376
ers, reducing interference between the forget and 4377
retention objectives; (b) the unlearning runs for 4378
40 steps (vs. 30 here) with a tuned $\lambda_a \in [40, 160]$; 4379
and (c) the retain set is curated *structural neighbors*, 4380
not general knowledge, avoiding direct competition 4381
with the target gradient. 4382

Result 4: GA confirms Remark 3. GA without 4383
AR (even at reduced $\eta = 2 \times 10^{-5}$) enters Regime B: 4384
 $p_{\text{max}} = 1.0$, entropy = 0, generations degenerate to 4385
single-token repetition (‘‘osos...’’). GA+AR pro- 4386
duces Regime A: the model confidently generates 4387
wrong tokens ($p_{\text{max}} = 0.993$, entropy = 0.06) but 4388
retains structural coherence (e.g., ‘‘The capital of 4389
Z capital of Z...’’). Neither GA variant approaches 4390
uniform: the spectral norm for GA+AR (0.0091) 4391
is three orders of magnitude above the theoretical 4392
uniform value (6.58×10^{-6}), confirming that the 4393
parametric trace of learned response structure is 4394
preserved under GA. 4395

The entropy–spectral-norm diagnostic. These 4396
four conditions validate the two-dimensional diag- 4397
nostic ($H, \|\Lambda\|_{\text{op}}$) proposed in Remark 3: 4398

- **KL-only:** (11.93, 1.03×10^{-5}) — near- 4399
uniform, Fisher-suppressed; 4400
- **KL+AR:** (0.28, 0.046) — target-channel uni- 4401
formization is strongly slowed by AR, while 4402
general knowledge is preserved; 4403
- **GA+AR:** (0.06, 0.009) — concentrated on 4404
wrong tokens, parametric trace intact; 4405
- **GA-only:** (0.00, 10^{-22}) — degenerate col- 4406
lapse, model destroyed. 4407

A single scalar (spectral norm or entropy alone) 4408
conflates genuinely uninformative outputs with de- 4409
generate or wrong-answer outputs; the pair uniquely 4410
identifies each regime. 4411

As a lightweight robustness check, a separate 4412
restricted-update probe on Qwen2.5-0.5B (last- 4413
block-only updates) reproduced the same qualita- 4414
tive ordering: KL-to-uniform raised target entropy 4415
from 6.47 to ~ 10.52 bits while reducing $\|\Lambda\|_{\text{op}}$ 4416
from 0.142 to ~ 0.0506 , whereas GA and GA+AR 4417
remained in a lower-entropy regime (~ 5.15 bits) 4418
with larger target covariance norms (~ 0.249). We 4419
do not treat this lightweight probe as a replacement 4420
for Table 74, nor as new primary evidence for AR- 4421
based retention preservation; rather, it confirms that 4422

4423 the entropy–spectral-norm diagnostic remains qual-
 4424 itatively stable even under a much smaller restricted-
 4425 update setting.

4426 **A note on parameter-level Fisher.** In v1 of this
 4427 experiment, the *parameter-level* diagonal Fisher
 4428 (empirical Fisher, computed as $\sum_i (\nabla_{\theta} \ell_i)^2$) was
 4429 also measured and showed the *opposite* trend: it
 4430 increased by 39× after KL-to-uniform. This is ex-
 4431 pected and does not contradict Theorem 3. Theo-
 4432 rem 3 concerns the *output-level* Fisher $\mathcal{F} = J^T \Lambda J$,
 4433 which factors through the softmax covariance Λ .
 4434 The parameter-level empirical Fisher measures the
 4435 average squared gradient of the *loss function* (cross-
 4436 entropy), which increases when the output becomes
 4437 uniform because cross-entropy loss against the
 4438 (now-uniform) prediction is large. The distinction
 4439 is between “how much does the *output distribu-*
 4440 *tion* depend on θ ” (Theorem 3, controlled by Λ)
 4441 vs. “how steep is the *loss landscape*” (empirical
 4442 Fisher, which can increase even as Λ is suppressed).

4443 D.15 Component Ablation: Saturation Curve 4444 and Early Stopping

4445 To quantify each Q-ROU component’s contribution
 4446 and identify the optimal stopping point, we conduct
 4447 a cumulative ablation study on Qwen2.5-0.5B with
 4448 periodic checkpointing. Three fictional facts are
 4449 injected (80 steps, $\eta = 5 \times 10^{-5}$), then unlearned
 4450 under seven conditions: (i) baseline, (ii) KL-to-
 4451 uniform only (40 steps), (iii) KL+AR (40 steps),
 4452 (iv) KL+AR+SLUG with injected-model AR ref-
 4453 erence, (v) KL+AR+SLUG with pre-injection AR
 4454 reference, (vi) full Q-ROU with injected reference,
 4455 and (vii) full Q-ROU with pre-injection reference.
 4456 Conditions (iv)–(vii) run for 160 steps with metrics
 4457 recorded every 20 steps.

4458 **Key finding 1: SLUG unlocks selective forget-**
 4459 **ting.** KL-only and KL+AR both achieve target
 4460 $p_{\max} < 10^{-4}$ but *destroy all knowledge*: general
 4461 $\|\Lambda\|_{\text{op}}$ collapses from 0.178 to 10^{-4} . At step 20,
 4462 KL+AR+SLUG achieves target $p_{\max} = 5 \times 10^{-3}$
 4463 (< 0.01 threshold) while preserving general $\|\Lambda\|_{\text{op}}$
 4464 at 0.052 (29% of baseline). This confirms that
 4465 SLUG (selective layer freezing) is the critical mech-
 4466 anism for selective forgetting.

4467 **Key finding 2: early stopping is essential at**
 4468 **small scale.** While target p_{\max} improves mono-
 4469 tonically with more steps, general knowledge de-
 4470 grades rapidly after step 20: general $\|\Lambda\|_{\text{op}}$ drops
 4471 from 0.052 at step 20 to 2.2×10^{-4} at step 40—a

Table 75: Component ablation on Qwen2.5-0.5B ($|\mathcal{V}| = 151,936$). All SLUG (and full Q-ROU) conditions cross the $\tau = 0.01$ threshold by step 20, but continued training degrades general knowledge. Pre-injection vs. post-injection AR references yield indistinguishable results at this injection scale.

Method	Step	Target p_{\max}	$H/\ln \mathcal{V} $	General $\ \Lambda\ _{\text{op}}$
Post-inject (baseline)	0	0.9999	0.0001	0.178
KL-Uniform	40	10^{-4}	0.986	10^{-4}
KL+AR	40	10^{-4}	0.986	10^{-4}
KL+AR+SLUG	20	5.0×10^{-3}	0.894	0.052
	40	1.4×10^{-4}	0.981	2.2×10^{-4}
	80	8×10^{-5}	0.988	1.2×10^{-4}
	160	4×10^{-5}	0.993	7.6×10^{-5}
Full Q-ROU	20	4.7×10^{-3}	0.894	0.055
	40	1.4×10^{-4}	0.981	2.2×10^{-4}
	80	8×10^{-5}	0.988	1.2×10^{-4}
	160	4×10^{-5}	0.993	7.6×10^{-5}

237× collapse. This indicates that on 0.5B mod- 4472
 4473 els, the SLUG subspace does not provide sufficient
 4474 isolation to prevent AR from being overwhelmed
 4475 by the KL-uniform gradient over extended train-
 4476 ing. At 3B/7B scale (Table 77), WMDP accuracy
 4477 is perfectly preserved ($\Delta = 0$) after 20-step Q-ROU,
 4478 confirming that this early-stopping sensitivity is
 4479 scale-dependent.

4480 **Key finding 3: QuantNoise is compatible with**
 4481 **SLUG-based unlearning.** With properly cali-
 4482 brated noise scale ($\epsilon_q = |\bar{w}| \times 0.01$), the
 4483 full Q-ROU condition (KL+AR+SLUG+Quant-
 4484 Noise+EWC) matches the KL+AR+SLUG satura-
 4485 tion curve at every checkpoint. The QuantNoise per-
 4486 turbation at this relative scale ($\sim 1\%$ of per-tensor
 4487 weight magnitude) does not interfere with the forget-
 4488 ting signal while still providing the loss-landscape
 4489 flattening validated in Section D.16.

4490 D.16 Empirical Verification of 4491 Loss-Landscape Flattening (Theorem 1)

4492 Theorem 1 predicts that QuantNoise training en-
 4493 courages flat loss basins, reducing sensitivity to
 4494 weight perturbations (including quantization noise).
 4495 We verify this by measuring *perturbed sharpness*,
 4496 defined as:

$$S = \frac{\mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [\mathcal{L}(\theta + \epsilon) - \mathcal{L}(\theta)]}{\sigma^2} \approx \frac{1}{2} \text{tr}(H) \quad (34)$$

4497 where H is the Hessian and the approximation holds
 4498 by Taylor expansion. This formulation avoids the
 4499 numerical instabilities of Hutchinson-based $\text{tr}(H)$
 4500

Table 76: Perturbed-sharpness audit for Theorem 1 on Qwen2.5-0.5B. QuantNoise reduces the estimated sharpness of the 30-step checkpoint relative to vanilla training under the same perturbation scale.

Condition	Sharpness S (mean \pm std)	Ratio vs. Untrained
Untrained	9,882 \pm 24,205	1.00
Vanilla (30 steps)	6,401 \pm 17,166	0.65
QuantNoise (30 steps)	3,810 \pm 17,426	0.39

4501 estimation in bfloat16 models.

4502 **Experimental setup.** We train the last two trans-
4503 former layers of Qwen2.5-0.5B for 30 steps on five
4504 factual statements, under two conditions: standard
4505 optimization (vanilla) and QuantNoise ($p = 0.1$).
4506 Sharpness is estimated using 30 Gaussian perturba-
4507 tions with $\sigma = 10^{-3}$.

4508 **Results.** QuantNoise training reduces perturbed
4509 sharpness by **40%** relative to vanilla training at
4510 the same step count ($S_{\text{QN}}/S_{\text{vanilla}} = 0.60$). This
4511 confirms that QuantNoise implicitly regularizes
4512 a noise-weighted Hessian trace, producing flatter
4513 loss basins that are more tolerant of the sampled
4514 post-training weight perturbations used in this au-
4515 dit, including INT4/NF4-like quantization noise.
4516 The large standard deviations reflect the inherent
4517 stochasticity of $O(10^7)$ -dimensional perturbation
4518 sampling with $n = 30$ draws; the *ordering* $S_{\text{QN}} <$
4519 $S_{\text{vanilla}} < S_{\text{untrained}}$ is consistent across all perturba-
4520 tion batches examined.

4521 D.17 MUSE/WMDP Utility Preservation

4522 To probe whether Q-ROU preserves general-
4523 purpose model capabilities, we evaluate
4524 general-knowledge accuracy (WMDP-style
4525 multiple-choice) and retain-domain perplexity
4526 before and after Q-ROU unlearning on Harry
4527 Potter-domain forget probes.

4528 **Key observations.** (1) WMDP accuracy is *un-*
4529 *changed* ($\Delta = 0.000$) on both 3B and 7B, con-
4530 sistent with low measured damage to this general-
4531 knowledge slice. (2) Retain PPL ratio decreases
4532 with scale (1.06 at 3B \rightarrow 1.02 at 7B), confirming
4533 the scale-dependent redundancy effect discussed in
4534 Section C.6: larger models have more capacity to
4535 absorb the forgetting gradient without disturbing
4536 retention. (3) Forget p_{max} on the 3B model reaches
4537 0.004 after 40 steps, below the strict $\tau = 0.01$
4538 threshold, indicating a target-side margin in this
4539 probe.

Table 77: Utility preservation after Q-ROU unlearning on two Qwen2.5 model scales. WMDP-style general-knowledge accuracy is perfectly preserved ($\Delta = 0$). Retain-domain perplexity increases by at most 6%, consistent with the neighbor-drift bounds predicted by Proposition 2.

Model	Condition	WMDP Acc	Forget p_{max}	Retain PPL Ratio
Qwen2.5-3B	Original	0.750	0.462	—
	Q-ROU	0.750	0.004	1.060
Qwen2.5-7B	Original	0.850	—	—
	Q-ROU	0.850	—	1.018

Relative KnowMem Reduction interpretation. 4540

4541 For MUSE, we report Relative KnowMem Reduc-
4542 tion ($1 - \text{KnowMem}_{\text{after}}/\text{KnowMem}_{\text{base}}$), so higher
4543 indicates stronger forgetting. Under this defini-
4544 tion, a retain/retrain baseline is expected to be near
4545 zero, with small nonzero offsets attributable to
4546 finite-sample ROUGE variability. The 3B Q-ROU
4547 value (0.0060) is therefore consistent with a mini-
4548 mal reduction at 40 steps, which may reflect under-
4549 saturation at this early checkpoint; we did not run a
4550 3B step-sweep due to budget. By contrast, the 8B
4551 160-step run yields a clearly nontrivial reduction
4552 (0.3540), indicating that the pipeline can meaning-
4553 fully move the KnowMem reduction metric when
4554 trained longer.

4555 E VulnScore: Quantifying Structural 4556 Interference

4557 To quantitatively support the distinction between
4558 semantic and structural neighbors discussed in Sec-
4559 tion D, we developed the *VulnScore* metric to pre-
4560 dict the interference magnitude on seemingly unre-
4561 lated concepts during unlearning.

4562 **Theoretical Motivation.** When unlearning a tar-
4563 get concept T via gradient updates $\Delta\theta_T$, the inter-
4564 ference experienced by a disjoint concept C can be
4565 approximated using the Fisher Information Matrix
4566 (FIM) of C . If F_C is the diagonal approximation
4567 of the empirical Fisher matrix computed on con-
4568 cept C 's data, the raw vulnerability can be defined
4569 as $VS_{\text{raw}}(C|T) = \Delta\theta_T^\top F_C \Delta\theta_T$. This approximates
4570 the KL-divergence (or second-order log-likelihood
4571 damage) inflicted on C by the update $\Delta\theta_T$. However,
4572 raw VS_{raw} is heavily biased by globally fragile con-
4573 cepts. To isolate the *specific entanglement* between
4574 T and C , we define the Specific Entanglement Ratio

4575 (SER):

$$4576 \quad \text{SER}(C|T) = \frac{\Delta\theta_T^\top F_C \Delta\theta_T}{\Delta\theta_C^\top F_C \Delta\theta_C} \quad (35)$$

4577 This normalizes the interference by the vulnerabil-
4578 ity of C to its own target update.

4579 **Methodology and Absolute Magnitude.** We
4580 evaluated a symmetric grid of 50 diverse concepts
4581 across 6 domains (Harry Potter, Lord of the Rings,
4582 Star Wars, Marvel, General Geography, etc.) on
4583 Llama-3.2-3B. Crucially, when measuring the em-
4584 pirical damage (the actual change in log-probability
4585 $\Delta \log P_C$), we discovered that aligned gradients can
4586 cause the neighbor’s probability to *improve* (a syn-
4587 ergistic effect). To construct a strictly linear corre-
4588 lation metric, we must measure the absolute magni-
4589 tude of structural disruption relative to the target’s
4590 improvement, defined as the Absolute Relative Em-
4591 pirical Damage (Abs RED):

$$4592 \quad \text{Abs RED}(T \rightarrow C) = \frac{|\Delta \log P_C|}{|\Delta \log P_T|} \quad (36)$$

4593 To ensure the empirical measurements remain
4594 within the valid linear-approximation regime of
4595 the Taylor expansion (preventing severe out-
4596 put collapse), we conducted an automated grid
4597 search over micro-perturbation scales ($\eta \in$
4598 $\{1e-3, 5e-4, 1e-4\}$).

4599 **Results: Near-Perfect Rank Correlation.** At
4600 the optimal perturbation scale ($\eta = 0.001$), the
4601 Spearman rank correlation (ρ) between the theo-
4602 retical prediction (SER) and the actual parameter
4603 disruption (Abs RED) across the 50×50 evalua-
4604 tion grid (2,450 pairs) was an exceptionally high
4605 **0.921** ($p < 10^{-60}$). This provides strong evidence
4606 that VulnScore is a highly accurate mathematical
4607 predictor for structural interference.

4608 **Structural Hubs vs. Semantic Distance.** The
4609 analysis revealed specific, robust ‘‘structural hubs’’.
4610 For instance, updates to HP_Voldemort or
4611 HP_Snape reliably inflicted substantial, asym-
4612 metric damage on entirely unrelated concepts like
4613 Gen_Geography or Marvel_Spiderman. When
4614 comparing the SER against the cosine similarity
4615 of the sentence embeddings for the concepts, the
4616 correlation was exactly $\rho \approx 0$. This empirically
4617 supports the central premise of Section D: *semantic*
4618 *isolation does not guarantee structural safety*. Con-
4619 cepts that are far apart in human-defined semantic
4620 space can still collide strongly in the underlying
4621 parameter space of Llama-3.2-3B.

Implications for Q-ROU. This finding under- 4622
scores the inadequacy of relying solely on generic 4623
pre-training data or semantic similarity distances 4624
to construct retention penalties. Q-ROU’s Active 4625
Retention directly minimizes this predicted Fisher 4626
damage: even when the provided neighbor set is se- 4627
lected based on semantic human intuitions, anchor- 4628
ing their output probabilities constrains the update 4629
 $\Delta\theta_T$ toward directions that are locally low-impact 4630
under the aggregated neighbor Fisher metric (as 4631
shown in Proposition 2). The $\rho = 0.921$ discovery 4632
provides strong empirical backing for this geomet- 4633
ric protection strategy. 4634

4635 **E.1 Theoretical Integration: The Divergence** 4636 **of Semantic and Structural** 4637 **Neighborhoods**

4638 Traditional machine unlearning benchmarks have
4639 implicitly operated under the assumption of an iso-
4640 morphism between the semantic embedding space
4641 (how humans categorize concepts) and the neu-
4642 ral parameter space (how the model physicalizes
4643 them). However, recent investigations into the inter-
4644 nal structure of Large Language Models reveal that
4645 this assumption is fundamentally flawed, leading to
4646 severe blind spots in evaluating unlearning safety
4647 and efficacy.

4648 **The Geometric Structure of Concepts and Fea-**
4649 **ture Crystals.** Recent dictionary-learning and
4650 sparse-autoencoder analyses suggest that concepts
4651 in LLMs are not stored as isolated vectors but as
4652 distributed feature structure with non-trivial over-
4653 lap (Li et al., 2025). In that setting, concepts that
4654 appear semantically distinct can still share nearby
4655 parametric support through broader latent vari-
4656 ables. Consequently, optimization trajectories that
4657 only push down target-token probabilities can also
4658 damage neighboring structure, even when ordinary
4659 semantic-similarity checks do not reveal the prob-
4660 lem immediately.

4661 **Parametric Knowledge Traces vs. Behavioral**
4662 **Suppression.** This kind of overlap helps ex-
4663 plain why standard gradient-based unlearning can
4664 achieve behavioral suppression without stable con-
4665 ceptual separation. As demonstrated by Hong *et*
4666 *al.* (Hong et al., 2025), target probabilities may fall
4667 while the underlying *parametric knowledge traces*
4668 remain largely intact or only weakly shifted. That
4669 leaves targeted knowledge vulnerable to adversar-
4670 ial recovery and increases the risk that aggressive
4671 parameter shifts overwrite neighboring structure.

4672 Robust conceptual suppression therefore requires
4673 interventions that affect the target trace while ex-
4674 plicitly protecting nearby traces, rather than relying
4675 only on output masking. The concept of knowl-
4676 edge neurons (Dai et al., 2022), which identifies
4677 specific FFN neurons responsible for factual re-
4678 call, provides direct evidence that different facts
4679 can share overlapping neuronal substrates—a phe-
4680 nomenon that manifests as structural interference
4681 in our VulnScore analysis. This shared substrate
4682 explains why semantically distant concepts can ex-
4683 hibit high structural vulnerability.

4684 **Depth-wise Activation Topology.** The loca-
4685 tion of these parametric traces is heavily depth-
4686 dependent. Lv et al. (Wang et al., 2025a) math-
4687 ematically map the activation patterns of LLM pa-
4688 rameters, proving that shallow layers exhibit dense
4689 activation patterns responsible for domain-agnostic,
4690 general linguistic capabilities. In contrast, middle-
4691 to-deep layers exhibit highly sparse activation pat-
4692 terns, acting as specialized storage for domain-
4693 specific facts and distinct entities. Applying broad
4694 parameter updates across all layers inevitably de-
4695 grades the foundational dense representations in
4696 shallow layers.

4697 **Synthesizing the Q-ROU Theoretical Frame-
4698 work.** Viewed through this advanced mechanistic
4699 lens, the empirical design choices in Q-ROU repre-
4700 sent a direct alignment with the geometric reality of
4701 the parameter space across three key dimensions.

4702 First, **SLUG explicitly targets specific sparse
4703 traces** by restricting updates to highly selective
4704 middle-to-deep layers (e.g., layer subsets around
4705 17–26 in Llama-3.2-3B). This allows the mecha-
4706 nism to bypass the dense, domain-agnostic shallow
4707 layers entirely (Wang et al., 2025a), thereby prevent-
4708 ing general capability degradation while operating
4709 precisely where the targeted parametric knowledge
4710 traces reside.

4711 Second, **Active Retention explicitly protects
4712 nearby structure** by anchoring neighbor distribu-
4713 tions during the edit. As established in Proposi-
4714 tion 2, this biases the unlearning update $\Delta\theta_T$ to-
4715 ward directions that are approximately null under
4716 the neighbor-aggregated Fisher Information Matrix.
4717 Even when starting from human-provided semantic
4718 neighbor sets, the mechanism translates semantic
4719 anchoring into a local structural constraint that pro-
4720 tects shared functional substructure from perturba-
4721 tion.

4722 Third, **the Wronskian serves as a supplement-
4723 ary collision alarm** to complement the main an-
4724 alytical results. As an empirical proxy (Observa-
4725 tion 1, Lemma 1), the Wronskian risk signal ρ_k
4726 monitors structural interference online. When stan-
4727 dard gradient ascent begins to overwrite shared re-
4728 tained structure, the AR surrogate models show
4729 converging roots, causing $\rho_k \rightarrow 0$. The adaptive
4730 variant of Q-ROU uses this signal to halt the update
4731 trajectory before retained traces drift further (Hong
4732 et al., 2025).

4733 Ultimately, the discovery regarding structural
4734 interference—such as our VulnScore metric achiev-
4735 ing a strong rank correlation of $\rho = 0.921$ with
4736 parameter damage despite having virtually zero
4737 correlation with semantic distance—supports this
4738 paradigm. Q-ROU succeeds not merely because of
4739 scalar regularization, but because it mechanistically
4740 recognizes and restricts its optimization within the
4741 geometric layout of parametric knowledge traces.

4742 E.2 Appendix Summary

4743 This appendix provides extended diagnostics, ro-
4744 bustness analyses, and theoretical addenda that sup-
4745 port the primary claims of the main text. We first
4746 present detailed step-sweeps and sensitivity studies
4747 on the expanded 65-probe set, including compre-
4748 hensive ablations of the adaptive mechanisms and
4749 hyperparameter dependencies (Sections C.4, C.6;
4750 Tables 36, 40). The follow-up audit section adds
4751 QR-LoRA-style low-rank baselines and benign re-
4752 learning stress tests, clarifying that quantization-
4753 stable target forgetting, semantic-neighbor preser-
4754 vation, and recurrence resistance are separable eval-
4755 uation axes (Section C.3). The technical stability of
4756 the framework is then established through thresh-
4757 old sensitivity analyses and loss-landscape sharp-
4758 ness diagnostics, complemented by validation on
4759 hardware NF4 (Sections C.12, D.16, C.17). To
4760 isolate the underlying mechanisms, we conduct
4761 AR transplant controls, SLUG-layer ablations, and
4762 sequential-forgetting persistence tests (Sections C.5,
4763 C.16). Furthermore, we include external sanity
4764 checks using the MUSE and WMDP benchmarks
4765 alongside expanded domain-transfer diagnostics
4766 (Sections B.5, D.17). Finally, we provide detailed
4767 audits of baseline failure modes, computational
4768 overhead accounting, and formal mathematical ad-
4769 denda regarding Fisher/Wronskian analysis and
4770 reproducibility (Sections D.4, D.3, D.2, E, D.8,
4771 B). The depth probing and embedding analyses ar-
4772 guably represent the most significant mechanistic

4773 findings. The 100% conceptual suppression across
4774 reverse association, multistep reasoning, and in-
4775 context extraction—combined with marked target-
4776 specific divergence in hidden representations and
4777 resilience to adversarial CoT—shows that Q-ROU
4778 goes beyond surface-level probability masking in
4779 the tested audit suite. Coupled with its computa-
4780 tional efficiency and quantization resilience, these
4781 properties establish Q-ROU as a viable framework
4782 for bounded post-deployment knowledge manage-
4783 ment while leaving robustness to arbitrary retrain-
4784 ing and stronger reconstruction attacks as open prob-
4785 lems.