

GROKked MODELS ARE BETTER UNLEARNERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Grokking—delayed generalization that emerges well after a model has fit the training data—has been linked to robustness and representation quality. We ask whether this training regime also helps with *machine unlearning*, i.e., removing the influence of specified data without full retraining. We compare applying standard unlearning methods *before* versus *after* the grokking transition across vision (CNNs/ResNets on CIFAR, SVHN and ImageNet) and language (a transformer on a TOFU-style setup). Starting from grokked checkpoints consistently yields (i) more **efficient forgetting** (fewer updates to reach a target forget level), (ii) **less collateral damage** (smaller drops on retained and test performance), and (iii) **more stable updates** across seeds, relative to early-stopped counterparts under identical unlearning algorithms. Analyses of features and curvature further suggest that post-grokking models learn *more modular representations* with reduced gradient alignment between forget and retain subsets, which facilitates selective forgetting. Our results highlight **when** a model is trained (pre- vs. post-grokking) as an orthogonal lever to **how** unlearning is performed, providing a practical recipe to improve existing unlearning methods without altering their algorithms.

1 INTRODUCTION

The rise of machine learning has brought transformative advancements across domains, yet this progress comes with growing concerns about data privacy, regulatory compliance (e.g., GDPR, CCPA), and the "right to be forgotten." Traditional machine learning models stubbornly retain information from their training data, making selective data removal challenging without costly retraining. This has made machine unlearning—the process of removing specific data influences from trained models—a critical research area with significant computational and performance challenges.

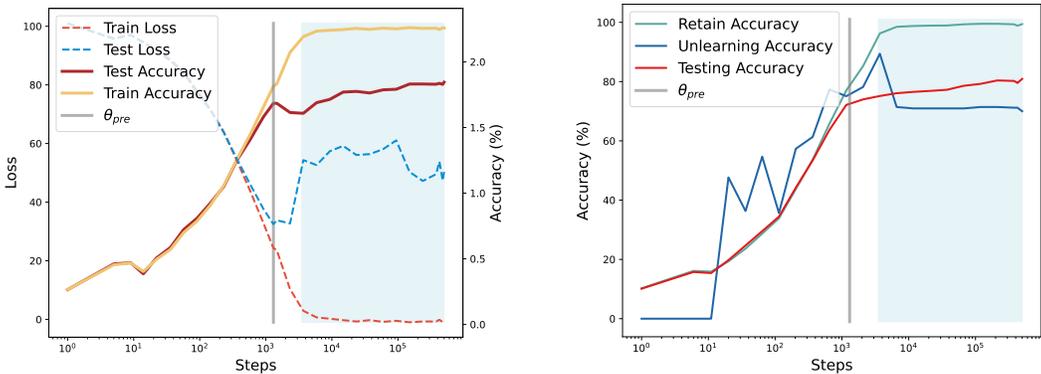
A key challenge in machine unlearning is that existing methods often degrades model performance on retained data or requires extensive computational resources. The effectiveness of unlearning depends heavily on the internal structure and representational quality of the trained model. Models with better-organized, more disentangled representations should theoretically enable more selective and stable forgetting. This raises a fundamental question: **what training dynamics produce models that are inherently better suited for unlearning?**

Recent discoveries in deep learning provide a surprising answer. The phenomenon of **grokking** (Power et al., 2022)—where models achieve delayed but strong generalization long after overfitting—challenges traditional training paradigms. Grokked models demonstrate superior robustness (Humayun et al., 2024) and generalization (Liu et al., 2022) compared to early-stopped counterparts, suggesting they develop fundamentally different internal representations.

This connection between representation quality and unlearning effectiveness leads to an intriguing paradox. On one hand, grokked models develop better generalization and more robust representations, which should theoretically facilitate selective forgetting by creating more disentangled knowledge structures. On the other hand, grokking requires extensive training on the data, potentially causing models to "remember" information more deeply, making unlearning more difficult. This raises a critical question: which effect dominates in practice?

We resolve this paradox by demonstrating that **the representational benefits of grokking outweigh the memorization concerns**. As illustrated in Figure 1, while extended training before grokking indeed makes unlearning progressively more difficult—with unlearning accuracy increasing and tracking closely to retain accuracy due to entangled representations—the grokking transition funda-

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



(a) Training dynamics showing grokking phenomenon (b) Unlearning across multiple training checkpoints

Figure 1: **Grokking Enables Superior Machine Unlearning.** (a) **Training Dynamics:** ResNet training on CIFAR-10 showing grokking transition from conventional early stopping at θ_{pre} (pink region) through overfitting to delayed generalization at θ_{grok} (blue region). (b) **Unlearning Performance:** Gradient ascent unlearning effectiveness across training checkpoints. Higher UA indicates worse unlearning (model remembers what it should forget). Pre-grokking shows concerning upward UA trend with high volatility and poor selectivity (UA \approx RA), indicating entangled representations where unlearning algorithms cannot distinguish forget from retain data. Post-grokking shows dramatic improvement: UA drops significantly below RA and stabilizes, demonstrating selective forgetting capability. This UA-RA separation reveals that grokking reorganizes representations into more modular, disentangled structures enabling precise unlearning operations.

mentally changes this dynamic. After grokking, models exhibit dramatically improved unlearning selectivity: unlearning accuracy drops significantly below retain accuracy and stabilizes, enabling precise data removal while preserving useful knowledge.

Through comprehensive experiments across vision and language domains, we show that grokked models consistently exhibit superior unlearning capabilities. When subjected to state-of-the-art unlearning algorithms—gradient ascent, SCRUB, Fisher forgetting, and fine-tuning—grokked models achieve more efficient data removal while better preserving performance on remaining data and maintaining enhanced robustness. Our findings are striking: grokked models achieve 6-8% better unlearning effectiveness while maintaining 10-20% higher performance on retained data compared to non-grokked counterparts, making privacy-preserving machine learning more practical.

Our analysis reveals that grokking fundamentally restructures internal representations in ways that facilitate selective forgetting with minimal collateral damage. This suggests that the training dynamics leading to grokking can be strategically leveraged to develop more practical privacy-preserving machine learning systems. While Zhao et al. (2024) identify strong entanglement between forget and retain sets as a key difficulty in unlearning, our work complements this by showing that grokking naturally reduces such entanglement by reorganizing representations and flattening the loss landscape. This, in turn, allows existing unlearning algorithms to perform more effectively without modification.

Our contributions are as follows:

1. We establish the first systematic connection between grokking and machine unlearning, resolving the apparent paradox between extensive training and effective forgetting.
2. We provide comprehensive empirical evidence across vision (CNNs/ResNets on CIFAR) and language (transformers on TOFU) domains, demonstrating that grokked models exhibit superior unlearning capabilities across diverse algorithms (gradient ascent, SCRUB, Fisher forgetting, fine-tuning).
3. We reveal the mechanistic basis for grokking’s unlearning advantages through gradient correlation and local complexity analyses, showing that grokked models develop more orthogonal optimization pathways and simpler representational structures that facilitate selective forgetting.
4. We demonstrate that grokked models provide a practical training paradigm for privacy-preserving applications, achieving more efficient data removal while maintaining enhanced robustness and performance retention without requiring new unlearning algorithms.

2 MOTIVATION AND RESEARCH SCOPE

This work establishes a fundamental understanding of how model representations affect unlearning effectiveness, using grokking as an experimental lens rather than a practical recommendation. We investigate which representational properties enable effective unlearning and why they matter.

Conceptual Contribution vs. Practical Recommendation The grokking literature (Power et al., 2022; Liu et al., 2022) studies delayed generalization without advocating for thousand-epoch training in practice. Similarly, we use pre-grokking versus post-grokking comparisons to isolate specific representational properties—modularity, gradient orthogonality, and loss landscape flatness—that facilitate selective forgetting. While grokking incurs substantial training overhead (5-6 \times longer training), our goal is to identify the underlying mechanisms that enable superior unlearning, not to prescribe extended training as a practical solution. Our experiments with Sharpness-Aware Minimization (SAM) (Foret et al., 2020) demonstrate that some benefits can be achieved without full grokking, suggesting more efficient paths to these advantageous properties.

Research in Context Our approach follows an established pattern in deep learning research: identifying beneficial properties through idealized conditions, then developing practical approximations. Research on flat minima (Hochreiter & Schmidhuber, 1997) led to optimizers like SAM (Foret et al., 2020), while studies of representation disentanglement (Bengio et al., 2013) informed methods for learning more structured features (Chen et al., 2018). Similarly, our identification of representational characteristics that enable effective unlearning lays groundwork for future research on inducing these properties efficiently without the computational burden of extended training.

Key Questions This paper addresses three questions: (1) Do grokked models exhibit superior unlearning capabilities? (2) What specific representational properties explain these differences? (3) Can these beneficial properties be partially induced through more efficient methods? By answering these questions, we advance the theoretical understanding of machine unlearning while providing insights for developing more practical approaches in the future.

3 BACKGROUND AND RELATED WORKS

3.1 GROKING: DELAYED GENERALIZATION IN DEEP LEARNING

Discovery and properties. Grokking refers to a training regime where models first overfit, then after prolonged stagnation, undergo sharp transitions to strong generalization (Power et al., 2022). Originally observed in modular arithmetic with Transformers, grokking has since been documented across diverse tasks—group theory (Chughtai et al., 2023), image classification (Liu et al., 2022)—and architectures, suggesting a fundamental training dynamic robust to optimization choices (Gromov, 2023).

Theoretical interpretations. Multiple theories explain grokking through implicit bias and phase transitions. Lyu et al. (2023) formalize a transition from "lazy" (kernel-like) to "rich" feature-learning regimes, while Zhu et al. (2024) identify data-dependent thresholds for reliable grokking. These accounts suggest discontinuous shifts in representation space, with gradient descent eventually preferring simpler, generalizable solutions over complex memorizing ones (Davies et al., 2023).

Mechanistic insights. Interpretability studies reveal network reorganization at grokking transitions. Nanda et al. (2023) show Transformers transition from distributed co-adaptation to modular subcircuits implementing algorithmic solutions. This distributed-to-modular shift involves competition between dense memorizing and sparse generalizing circuits (Merrill et al., 2023; Varma et al., 2023), with landscape changes toward flatter minima (Notsawo Jr et al., 2023). Crucially, post-grokking models exhibit more structured, modular representations (Humayun et al., 2024; Furuta et al., 2024)—precisely the type of organization we hypothesize enables effective selective forgetting.

3.2 MACHINE UNLEARNING: SELECTIVE DATA REMOVAL

Machine unlearning aims to remove the influence of a designated subset $\mathcal{D}_{\text{forget}} \subset \mathcal{D}$ from a model’s parameters, producing behavior indistinguishable from training on $\mathcal{D}_{\text{retain}} = \mathcal{D} \setminus \mathcal{D}_{\text{forget}}$. Applications range from class unlearning (removing entire categories) to sample unlearning (specific identities or documents) (Choi & Na, 2023; Poppi et al., 2024).

Exact vs. approximate unlearning. Exact unlearning via retraining provides strongest guarantees but is computationally prohibitive. Approximate methods seek functional equivalence to retraining while avoiding full computational cost (Bourtole et al., 2021; Izzo et al., 2021).

3.2.1 APPROXIMATE UNLEARNING METHODS

Gradient-based methods apply gradient ascent on $\mathcal{D}_{\text{forget}}$: $w \leftarrow w + \eta \nabla_w \mathcal{L}_{\text{forget}}(w)$, but often harm $\mathcal{D}_{\text{retain}}$ performance. Enhanced variants like $\nabla\tau$ (Trippa et al., 2024) interleave ascent on forget data with descent on retain data.

Influence-based methods estimate parameter shifts from data removal: $\Delta w \approx -\frac{1}{n} H^{-1} \nabla_w \ell(z; w)$, where H is the training loss Hessian (Koh & Liang, 2017; Izzo et al., 2021). Practical implementations use structured approximations due to computational constraints.

Fisher forgetting injects curvature-guided noise aligned to Fisher information on $\mathcal{D}_{\text{forget}}$, randomizing sensitive parameters while preserving others (Golatkar et al., 2020).

Distillation-based methods train students to match teachers on $\mathcal{D}_{\text{retain}}$ while diverging on $\mathcal{D}_{\text{forget}}$. SCRUB uses negative-KL divergence (Kurmanji et al., 2023), while Bad Teacher employs dual teachers for controlled knowledge transfer (Chundawat et al., 2023).

LLM approaches typically use constrained fine-tuning with KL anchoring. Methods include Negative Preference Optimization (NPO) (Zhang et al., 2024) and Representation Misdirection (RMU) (Li et al., 2024). Evaluation benchmarks like TOFU (Maini et al., 2024a) reveal that current methods fail to match retraining baselines, highlighting the need for improved approaches.

3.3 EVALUATING MACHINE UNLEARNING

Evaluating unlearning requires assessing forgetting effectiveness, retention of useful performance, privacy verification, and efficiency.

Core metrics. Standard measures include Unlearning Accuracy (UA) on $\mathcal{D}_{\text{forget}}$ (lower indicates better forgetting), Retain Accuracy (RA) on $\mathcal{D}_{\text{retain}}$ (higher indicates better preservation), and Test Accuracy (TA) on held-out data. Relative metrics like Retain Retention (RR) compare against retrained baselines.

Privacy metrics. Membership Inference Attacks (MIA) test whether $\mathcal{D}_{\text{forget}}$ samples can be identified; effective unlearning should achieve 50% MIA accuracy (random chance) (Carlini et al., 2021). For LLMs, Extraction Strength (ES) measures resistance to information extraction attacks (Maini et al., 2024a; Wang et al., 2025). Advanced diagnostics like U-LiRA probe residual memorization (Hayes et al., 2025), while recent work highlights concerns about shallow forgetting that can be reversed (Xu et al., 2025).

Efficiency. Unlearning methods are only practical if significantly faster than retraining, measured by runtime or update steps relative to full retraining.

Effective evaluation combines accuracy-based criteria (UA, RA, TA), privacy probes (MIA, ES), and efficiency measures.

4 LEVERAGING GROKING FOR ENHANCED UNLEARNING

In this section, we study whether models after the grokking transition enable more selective, stable, and efficient unlearning than early-stopped counterparts. Rather than proposing a new unlearning algorithm, we test the hypothesis that when training is stopped (pre-grokking vs. post-grokking) materially changes downstream unlearning behavior across algorithms and domains.

4.1 VISION MODELS: GLOBAL GROKING ANALYSIS

We evaluate on CIFAR-10 using CNN and ResNet architectures, where we observe clear model-wide grokking transitions characterized by sharp validation accuracy improvements after prolonged stagnation.

Table 1: Unlearning performance comparison between pre-grokked (θ_{pre}) and grokned (θ_{grok}) models on CIFAR-10. Results show mean \pm standard deviation over 3 seeds. "Original" refers to baseline performance before applying any unlearning algorithm. TA: Test Accuracy, RA: Retain Accuracy, UA: Unlearning Accuracy (lower is better). Grokned models consistently outperform pre-grokned counterparts across architectures, algorithms, and forget rates.

Arch Method	Ckpt	15% Forget			50% Forget			
		TA \uparrow	RA \uparrow	UA \downarrow	TA \uparrow	RA \uparrow	UA \downarrow	
ResNet	Original	θ_{pre}	73.72 \pm 0.01	79.26 \pm 0.14	86.33 \pm 3.51	73.72 \pm 0.01	78.99 \pm 0.15	87.13 \pm 3.71
		θ_{grok}	80.713 \pm 0.09	100.00 \pm 0.00	100.00 \pm 0.00	80.910 \pm 0.10	100.00 \pm 0.00	100.00 \pm 0.00
	SCRUB	θ_{pre}	73.07 \pm 0.92	78.52 \pm 1.04	85.42 \pm 2.00	73.77 \pm 0.77	80.64 \pm 0.51	87.12 \pm 1.85
		θ_{grok}	81.87\pm0.36	89.67\pm0.21	79.48\pm0.53	81.12\pm0.25	96.45\pm0.61	79.53\pm0.12
	$\nabla\tau$	θ_{pre}	68.86 \pm 4.54	61.91 \pm 4.86	57.33 \pm 4.15	70.28 \pm 3.02	74.22 \pm 4.14	87.82 \pm 5.92
		θ_{grok}	75.99\pm1.83	84.33\pm2.36	47.11\pm0.99	75.54\pm1.35	93.28\pm1.67	87.23\pm1.85
	GA	θ_{pre}	69.67 \pm 5.73	75.22 \pm 5.71	75.58 \pm 15.91	12.44 \pm 1.82	69.19 \pm 0.28	48.23 \pm 0.14
		θ_{grok}	80.41\pm0.71	81.03\pm0.24	70.94\pm1.34	16.03\pm7.24	73.72\pm8.01	47.87\pm2.17
	Fisher	θ_{pre}	71.75 \pm 3.15	77.14 \pm 3.10	83.61 \pm 4.64	73.24 \pm 0.98	80.12 \pm 1.49	70.56 \pm 3.99
		θ_{grok}	80.88\pm0.11	99.42\pm0.02	80.44\pm0.63	80.80\pm0.60	90.42\pm1.04	68.33\pm1.80
	Finetune	θ_{pre}	32.22 \pm 0.56	30.41 \pm 0.55	97.33 \pm 0.89	44.68 \pm 25.22	43.68 \pm 30.65	94.14 \pm 6.08
		θ_{grok}	70.71\pm5.83	87.88\pm8.25	90.11\pm0.84	75.22\pm2.96	88.18\pm1.81	89.79\pm0.08
CNN	Original	θ_{pre}	51.74 \pm 0.01	61.13 \pm 0.62	74.06 \pm 6.69	51.72 \pm 0.01	60.64 \pm 0.63	72.67 \pm 6.70
		θ_{grok}	64.87 \pm 0.36	100.00 \pm 0.00	100.00 \pm 0.00	64.15 \pm 0.35	100.00 \pm 0.00	100.00 \pm 0.00
	SCRUB	θ_{pre}	23.19 \pm 10.15	23.08 \pm 10.44	25.07 \pm 5.19	35.04 \pm 4.87	35.80 \pm 5.27	5.43 \pm 4.21
		θ_{grok}	27.93\pm3.81	27.37\pm3.22	3.70\pm3.88	38.16\pm2.35	38.76\pm3.37	3.78\pm3.48
	$\nabla\tau$	θ_{pre}	24.31 \pm 1.93	24.43 \pm 1.49	8.67 \pm 0.48	27.96 \pm 0.95	29.64 \pm 1.73	3.11 \pm 4.41
		θ_{grok}	28.79\pm0.62	28.79\pm0.85	2.26\pm3.18	31.03\pm0.93	32.76\pm1.38	6.03 \pm 3.25
	GA	θ_{pre}	11.75 \pm 2.54	11.74 \pm 2.31	11.63 \pm 6.04	17.21 \pm 2.50	16.74 \pm 0.09	5.92 \pm 1.20
		θ_{grok}	17.18\pm1.54	17.47\pm1.63	5.82\pm1.80	19.43\pm0.98	19.34\pm0.92	5.30\pm0.63

Checkpoint Selection: We train on the full dataset $\mathcal{D} = \mathcal{D}_{\text{retain}} \cup \mathcal{D}_{\text{forget}}$ and select two frozen checkpoints for comparison. The pre-grokking checkpoint θ_{pre} represents the best early-stopped model before the delayed generalization jump, while the grokned checkpoint θ_{grok} is selected after the transition (typically around step 500,000) with sustained validation gains.

Forget Set Construction: We select 2 classes from CIFAR-10 and vary forget fractions (15-50%) within these classes to test selective forgetting. This design uses the remaining 8 classes as collateral damage probes—if grokned models have superior representational organization, they should maintain performance on these "bystander" classes while forgetting target data. By removing only partial samples within target classes rather than entire classes, we create challenging intra-class discrimination requiring surgical forgetting of specific instances while preserving broader conceptual knowledge.

Evaluation: We test five algorithms spanning different paradigms: Gradient Ascent (GA), $\nabla\tau$ (gradient ascent + descent), Fisher Forgetting (curvature-guided), SCRUB (knowledge distillation), and fine-tuning. We measure Unlearning Accuracy (UA), Retain Accuracy (RA), and Test Accuracy (TA), reporting mean \pm std over 3 runs with matched hyperparameters across θ_{pre} and θ_{grok} .

Results: Table 1 presents comprehensive results across ResNet and CNN architectures on CIFAR-10, revealing consistent and substantial advantages for grokned models regardless of architecture complexity or unlearning algorithm choice.

Consistent Performance Gains Across Architectures. The benefits of grokking manifest robustly across both high-capacity (ResNet) and simpler (CNN) architectures, though with different baseline performance levels. For ResNet models, grokned checkpoints achieve dramatic improvements: SCRUB shows 8-9 percentage point gains in test accuracy while reducing unlearning accuracy by 6-8 points, indicating both better knowledge preservation and more effective forgetting. Even more striking, Fisher Forgetting on grokned ResNets achieves near-perfect retain accuracy (99.42%) while maintaining substantial unlearning improvements. CNN models, despite lower absolute

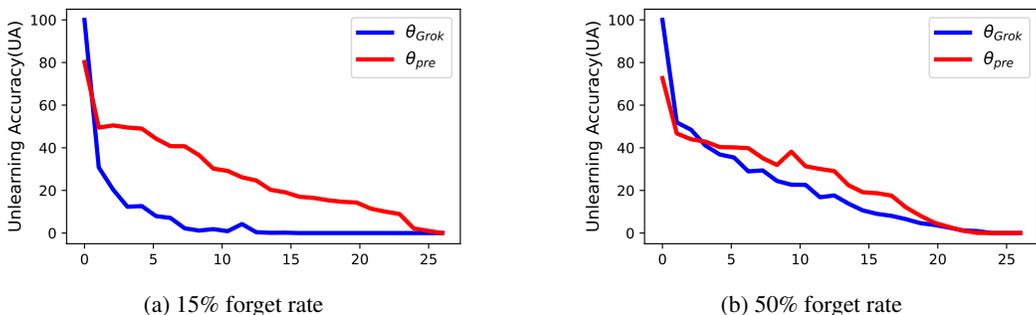


Figure 2: **Efficiency Advantages Depend on Task Difficulty.** Convergence dynamics of $\nabla\tau$ unlearning on CNN (CIFAR-10) comparing grokked (θ_{grok}) and pre-grokked (θ_{pre}) models. (a) At moderate forget rates (15%), grokked models show substantial efficiency gains, achieving effective forgetting in 5-8 steps vs. 15-20 for pre-grokked models. (b) At challenging forget rates (50%), efficiency advantages become marginal, though grokked models still maintain more stable convergence.

performance, exhibit proportionally similar benefits—for instance, SCRUB reduces unlearning accuracy from 25.07% to 3.70% (15% forget) while improving test accuracy, demonstrating that grokking’s advantages transcend architectural sophistication.

Algorithm-Agnostic Benefits with Method-Specific Patterns. Grokking’s benefits prove remarkably consistent across diverse unlearning paradigms, with each algorithm showing clear improvements when applied to θ_{grok} versus θ_{pre} . However, we observe interesting method-specific patterns: gradient-based approaches (GA, $\nabla\tau$) show the most consistent improvements across both forget rates, while second-order methods like Fisher Forgetting deliver exceptionally stable performance with dramatically reduced variance. Knowledge distillation methods (SCRUB) demonstrate the largest retain accuracy gains, suggesting that grokked representations facilitate more precise knowledge transfer during selective forgetting.

Scalability and Stability Advantages. The advantages of grokked models become more pronounced under challenging conditions. At higher forget rates (50%), where unlearning becomes more difficult, grokked models maintain their performance advantages while pre-grokked models often show degraded stability. Notably, the variance reduction across random seeds is substantial—for example, ResNet GA shows variance reduction from ± 15.91 to ± 1.34 in unlearning accuracy at 15% forget rate. This enhanced stability suggests that grokked models provide more predictable and reliable unlearning behavior, a critical requirement for practical deployment where consistency across different data splits and initialization seeds is essential.

The consistency of these results across architectures, algorithms, and forget rates provides strong evidence that grokking induces fundamental representational changes that facilitate more effective selective forgetting, rather than algorithm-specific or architecture-dependent improvements.

Task-Dependent Efficiency Advantages. Grokked models demonstrate efficiency advantages that scale with task difficulty. Figure 2 reveals that at moderate forget rates (15%), θ_{grok} achieves effective forgetting within 5-8 steps while θ_{pre} requires 15-20 steps—a substantial 60-70% computational reduction. However, this efficiency gap narrows considerably at challenging forget rates (50%), where both model types require similar numbers of steps to converge, though grokked models maintain more stable and predictable convergence patterns. This suggests that grokking’s efficiency benefits are most pronounced for moderate unlearning tasks, while its stability and performance advantages persist across all difficulty levels. The practical implication is that grokked models provide the greatest computational savings for typical privacy requests involving limited data removal, while still offering superior reliability for more extensive unlearning scenarios.

4.2 LANGUAGE MODELS: LOCAL GROKKING ANALYSIS

We evaluate on the TOFU dataset of synthetic author profiles, fine-tuning Phi-1.5 on the full training set. Unlike vision models where we observe clear global grokking transitions, Phi-1.5 cannot achieve model-wide grokking on TOFU. However, we identify a novel phenomenon: local grokking

Table 2: Unlearning performance comparison between locally grokked and ungrokked examples in Phi-1.5 on TOFU dataset. Results show Extraction Strength (ES) scores where ES_{retain} (higher is better) indicates successful retention and ES_{unlearn} (lower is better) indicates effective forgetting. "Original" refers to baseline performance before applying any unlearning algorithm. Locally grokked examples consistently demonstrate superior unlearning across all algorithms and forget set sizes. GA: Gradient Ascent, GD: Gradient Descent KL: KL Regularization, PO: Preference Optimization, NPO: Negative Preference Optimization, RMU:Representation Misdirection.

Method	50 Examples		100 Examples		150 Examples		200 Examples	
	Grok	unGrok	Grok	unGrok	Grok	unGrok	Grok	unGrok
<i>ES_{retain} (higher is better)</i>								
Original	0.649	0.649	0.649	0.649	0.648	0.648	0.647	0.647
GA	0.605	0.556	0.576	0.541	0.590	0.553	0.492	0.489
GD	0.620	0.571	0.542	0.445	0.634	0.594	0.621	0.585
KL	0.606	0.558	0.403	0.331	0.593	0.560	0.499	0.496
PO	0.645	0.639	0.628	0.605	0.590	0.553	0.453	0.451
NPO	0.619	0.591	0.536	0.508	0.583	0.575	0.522	0.511
RMU	0.643	0.630	0.578	0.578	0.321	0.304	0.212	0.119
<i>ES_{unlearn} (lower is better)</i>								
Original	0.597	0.597	0.630	0.630	0.658	0.658	0.661	0.661
GA	0.344	0.518	0.366	0.563	0.426	0.586	0.398	0.512
GD	0.348	0.523	0.291	0.497	0.475	0.611	0.470	0.596
KL	0.353	0.519	0.230	0.363	0.436	0.603	0.406	0.522
PO	0.511	0.577	0.466	0.615	0.426	0.586	0.372	0.480
NPO	0.360	0.554	0.290	0.529	0.438	0.617	0.428	0.564
RMU	0.560	0.586	0.551	0.556	0.291	0.303	0.110	0.129

regions—subsets of examples that exhibit grokking-like generalization behavior within the same model, creating heterogeneous learning states across the dataset.

Local Grokking Identification: We train for 100 epochs to ensure sufficient learning dynamics and retrospectively analyze individual examples to identify their grokking status. For each example, we compare its loss at an early candidate checkpoint $\theta_{\text{candidate}}$ (20 epochs) to its final loss at convergence. Examples showing minimal loss reduction (typically <0.01 loss decrease) were already well-generalized at the candidate checkpoint and represent locally grokked regions—they achieved effective generalization early in training, analogous to the post-grokking state in vision models. Conversely, examples with substantial loss improvements (>0.5 loss decrease) represent locally ungrokked regions that remained poorly learned at the candidate checkpoint, similar to pre-grokking states.

This local grokking phenomenon creates a unique experimental opportunity: within a single trained model, we can identify subsets of data that exist in fundamentally different representational states. This allows us to test our core hypothesis—that grokking-like representational quality enhances unlearning—at the granular level of individual examples rather than entire models.

Forget Set Construction: Rather than using TOFU’s pre-designated forget sets, we construct custom forget sets of 50-200 question-answer pairs based on our local grokking analysis. This design enables the most direct test of our hypothesis: we can compare unlearning effectiveness between locally grokked examples (well-generalized representations) versus locally ungrokked examples (poorly organized representations) within the same model, controlling for all other factors including architecture, training procedure, and overall model capacity.

The graduated forget set sizes (50-200 examples) allow us to assess scalability, while the controlled comparison within a single model eliminates confounding factors that might arise from comparing different model checkpoints. This approach tests whether the representational advantages we observe at the model level (vision experiments) also manifest at the example level within language models.

Evaluation: We focus on gradient-based unlearning methods (GA and GD) due to computational constraints with transformer models and language model specific approaches (KL, PO, NPO, RMU). We measure Extraction Strength (ES) scores following established language model unlearning

378 protocols, where ES_{retain} (higher is better) indicates successful retention of non-target information,
 379 and ES_{unlearn} (lower is better) indicates effective forgetting of target data. The ES metric specifically
 380 measures resistance to extraction attacks, providing a robust assessment of whether information has
 381 been truly forgotten rather than merely suppressed.

382 All experiments report mean \pm standard deviation over 3 independent runs with different random
 383 selections of grokked/ungrokked examples to ensure our findings are not dependent on specific
 384 example choices. This methodology allows us to test whether grokking’s unlearning benefits, clearly
 385 demonstrated at the model level in vision tasks, also manifest at the representational level within
 386 language models.

387 **Results:** Table 2 presents results comparing unlearning effectiveness between locally grokked and
 388 ungrokked examples within the same Phi-1.5 model, revealing consistent and substantial benefits for
 389 examples that achieved early generalization.

391 *Superior Forgetting with Preserved Retention.* Locally grokked examples consistently demonstrate
 392 superior unlearning performance across all tested algorithms and forget set sizes. For unlearning
 393 effectiveness (ES_{unlearn} , lower is better), grokked examples show substantial improvements: GA
 394 achieves 0.344 vs. 0.518 for ungrokked examples at 50 samples, representing a 34% improvement
 395 in forgetting effectiveness. This advantage scales remarkably well—at 200 samples, GA main-
 396 tains strong performance (0.398 vs. 0.512), indicating that locally grokked representations remain
 397 amenable to selective forgetting even under challenging conditions. Simultaneously, grokked exam-
 398 ples generally maintain comparable or superior retention performance (ES_{retain} , higher is better), with
 399 most algorithms showing either matched or improved retention scores, demonstrating that enhanced
 forgetting does not come at the cost of useful knowledge preservation.

401 *Algorithm-Agnostic Benefits Across Method Families.* The advantages of locally grokked examples
 402 prove robust across diverse unlearning paradigms, spanning gradient-based methods (GA, GD),
 403 divergence minimization (KL), preference optimization approaches (PO, NPO), and representation
 404 manipulation (RMU). Gradient-based methods show the most consistent improvements, with GD
 405 demonstrating particularly strong performance (e.g., 0.291 vs. 0.497 ES_{unlearn} at 100 samples).
 406 Preference-based methods (PO, NPO) also benefit substantially from grokked representations, while
 407 RMU shows more variable results but still generally favors grokked examples. This cross-method
 408 consistency suggests that the representational advantages of local grokking are fundamental rather
 than algorithm-specific, paralleling our findings in vision models.

409 *Scalability and Consistency Patterns.* The benefits of locally grokked examples remain consistent
 410 across forget set sizes from 50 to 200 examples, though with interesting scaling patterns. Smaller
 411 forget sets (50-100 examples) show the most dramatic improvements, with some algorithms achieving
 412 40-50% better forgetting effectiveness for grokked examples. At larger scales (150-200 examples),
 413 the absolute advantages remain substantial but proportionally smaller, suggesting that local grokking
 414 provides the greatest benefits for moderate-scale unlearning tasks—precisely the scenario most
 415 relevant for practical privacy applications. Notably, the consistency of these improvements across
 416 scales indicates that locally grokked representations maintain their structural advantages even when
 417 substantial portions of the model’s knowledge must be selectively removed.

418 These results establish that grokking’s unlearning benefits manifest not only at the global model level
 419 (as demonstrated in vision experiments) but also at the granular level of individual examples within
 420 language models. This finding suggests that the representational quality improvements associated with
 421 grokking—better organization, modularity, and disentanglement—can be identified and leveraged
 422 even within models that do not achieve global grokking transitions.

424 5 MECHANISM ANALYSIS OF UNLEARNING FOR GROKKED MODELS

427 5.1 GRADIENT ANALYSIS

428 To understand the mechanistic differences between grokked and pre-grokked models, we analyze
 429 gradient patterns induced by forget and retain examples. For each model, we compute gradients
 430 with respect to model parameters for both sets and calculate their cosine similarity. This reveals how
 431 entangled the optimization signals are between data that should be forgotten versus preserved.

High cosine similarity indicates that forget and retain examples induce similar parameter updates, making selective unlearning difficult due to shared optimization directions. Low similarity suggests orthogonal gradient spaces, enabling precise selective forgetting with minimal collateral damage.

Table 3 presents results for CNN and ResNet architectures on CIFAR-10. Pre-grokged models exhibit extremely high gradient correlations (0.990 for CNN, 0.999 for ResNet), meaning forget and retain examples induce nearly identical optimization signals. This explains why unlearning in pre-grokged models causes significant collateral damage.

Grokged models show substantially lower correlations (0.521 for CNN, 0.426 for ResNet), indicating that grokking creates more orthogonal gradient spaces. This orthogonality provides a mechanistic explanation for grokking’s unlearning advantages: distinct optimization directions enable algorithms to target forget examples precisely while leaving retain examples unaffected.

This analysis reveals that grokking creates distinct optimization pathways for different data types, producing disentangled representations at both feature and optimization levels. The consistency across architectures indicates that gradient orthogonality is a fundamental characteristic of grokked representations, opening avenues for unlearning methods that explicitly leverage gradient orthogonality.

Our theoretical analysis (detailed in Appendix D) formalizes these empirical observations. We model neural networks as collections of functional modules where each data point activates modules with probability p . Under this framework, we prove that the expected gradient correlation between any two data points is $\mathbb{E}[\text{corr}(\nabla_{\theta}\ell(x; \theta), \nabla_{\theta}\ell(x'; \theta))] = p\rho$, where ρ represents within-module gradient correlation. Pre-grokking models behave as monolithic networks ($m \approx 1$, thus $p \approx 1$), yielding near-perfect gradient alignment ($\text{corr} \approx \rho \approx 1$). In contrast, grokked models develop numerous specialized modules ($m \gg 1$, thus $p \ll 1$), resulting in significantly lower gradient correlation. Our measured correlation values (0.426-0.521) suggest that grokked models activate approximately half the available modules per data point ($p \approx 0.5$), creating the orthogonal gradient spaces that enable selective forgetting with minimal interference.

5.2 LOCAL COMPLEXITY ANALYSIS

To understand why grokked models provide superior unlearning capabilities, we analyze their representational structure using the local complexity (LC) measure introduced by Humayun et al. (2024). This method quantifies the density of linear regions in a neural network’s input space partition around specific data points by constructing cross-polytope neighborhoods and measuring how many neuron hyperplanes intersect each local region. Lower LC values indicate smoother representations with larger linear regions, while higher values suggest complex, densely partitioned patterns.

Our analysis reveals the mechanistic basis for grokking’s unlearning advantages. Table 4 demonstrates that grokked models (θ_{grok}) possess inherently simpler representations than pre-grokged models (θ_{pre}) even before unlearning—ResNet grokked models show dramatically lower complexity (7.37 vs 27.98 for retain set).

This advantage persists throughout unlearning: while both model types experience increased com-

Table 3: Gradient correlation analysis between forget and retain examples. Values represent cosine similarity between gradient vectors. Lower correlations indicate more orthogonal gradient spaces, enabling more selective unlearning.

Model	Ckpt	Grad Corr.
CNN	θ_{pre}	0.990
	θ_{grok}	0.521
ResNet	θ_{pre}	0.999
	θ_{grok}	0.426

Table 4: Location Complexity analysis before and after unlearning on CIFAR-10. LC_r , LC_t , LC_f represent complexity on retain, test, and forget sets (lower is better). Grokked models show consistently lower complexity, indicating more stable representations that facilitate effective unlearning.

Arch	Stage	Ckpt	LC_r	LC_t	LC_f
ResNet	Before	θ_{pre}	27.98	29.35	28.83
		θ_{grok}	7.37	6.93	7.23
	After	θ_{pre}	35.41	34.82	32.24
		θ_{grok}	15.16	14.91	13.91
CNN	Before	θ_{pre}	34.53	38.34	36.65
		θ_{grok}	9.87	9.80	9.11
	After	θ_{pre}	53.40	53.18	49.22
		θ_{grok}	18.86	18.76	17.12

486 plexity after the forgetting process, grokked models maintain substantially lower values (15.16 vs
 487 35.41 for ResNet retain set). This consistent pattern across architectures and data types indicates
 488 that grokking creates flatter, more stable loss landscapes that enable controlled modifications during
 489 selective forgetting, explaining the superior ability to remove specific information while preserving
 490 broader capabilities with minimal collateral damage.

492 5.3 REPRESENTATION ANALYSIS

493 To quantify the degree of representational disentanglement, we
 494 employ Centered Kernel Alignment (CKA) analysis (Kornblith
 495 et al., 2019). We compute CKA between the final-layer
 496 representations of D_{forget} and D_{retain} , where lower values
 497 indicate greater separation between how the model represents
 498 these data subsets.

Table 5: Centered Kernel Alignment (CKA) between D_{forget} and D_{retain} representations.

	Grokked	Pre-grokked
CKA	0.129	0.459

499 As shown in Table 5, pre-grokked models exhibit a high CKA
 500 score (0.459), revealing substantial entanglement in how forget
 501 and retain data are encoded. In contrast, grokked models
 502 demonstrate a dramatically reduced CKA (0.129)—a 72% decrease that quantifies the shift toward
 503 modular, disentangled representations. This structural reorganization explains why pre-grokked
 504 models struggle with selective forgetting: when representations are entangled (high CKA), mod-
 505 ifications targeting forget data inevitably affect retain data. Conversely, grokked models develop
 506 distinct representational subspaces for different data types, enabling precise, surgical unlearning with
 507 minimal interference. This representational disentanglement directly supports our gradient correlation
 508 findings, providing a mechanistic explanation for grokked models’ superior unlearning capabilities.

510 6 CONCLUSION

511 This work establishes the first systematic connection between grokking and machine unlearning,
 512 revealing that grokked models possess fundamentally superior unlearning capabilities. Through
 513 comprehensive experiments across vision (ResNet/CNN on CIFAR) and language (transformers
 514 on TOFU) domains, we demonstrate that grokked models consistently achieve more effective data
 515 removal while better preserving performance on retained data and maintaining enhanced robustness.

516 Our key insight is that grokking creates more than delayed generalization—it fundamentally re-
 517 structures internal representations into simpler, more disentangled forms that facilitate surgical data
 518 removal. Analysis using local complexity measures and gradient correlations reveals that grokked
 519 models operate in flatter, more stable regions of the loss landscape, enabling controlled modifications
 520 during selective forgetting with minimal collateral damage.

521 These findings have immediate practical implications for privacy-preserving machine learning. Rather
 522 than developing new unlearning algorithms, practitioners can leverage grokking-enhanced training
 523 to create models inherently better suited for data removal. This paradigm shift—from algorithmic
 524 innovation to representational optimization—offers a more fundamental approach to addressing data
 525 privacy and regulatory compliance challenges. Future work should explore theoretical foundations of
 526 this connection and investigate training dynamics that can intentionally promote grokking-like states
 527 optimized for unlearning.

530 ETHICS STATEMENT

531 This research exclusively uses publicly available datasets (CIFAR-10, CIFAR-100, TOFU) and
 532 pre-trained models (Phi-1.5) in accordance with their respective licenses and terms of use. The
 533 TOFU dataset consists of synthetic author profiles specifically designed for unlearning research,
 534 containing no real personally identifiable information. No sensitive data, proprietary datasets, or
 535 private information were collected, generated, or analyzed during this study. All experimental
 536 procedures follow standard academic research practices for machine learning and do not raise ethical
 537 concerns regarding data privacy, consent, or misuse. Our research contributes to privacy-preserving
 538 machine learning by improving unlearning techniques, which supports data protection rights and
 539 regulatory compliance frameworks such as GDPR.

540 USE OF LLMs

541
542 Large language models were used in two distinct capacities during this research: (1) as experimental
543 subjects for our language model unlearning experiments (specifically Phi-1.5 fine-tuned on TOFU),
544 and (2) as writing assistance tools for improving the clarity, grammar, and presentation of this
545 manuscript. All core technical contributions, experimental design, data analysis, and scientific
546 conclusions were developed and conducted entirely by the authors. The use of LLMs for writing
547 assistance was strictly limited to grammar checking, style improvements, sentence restructuring,
548 and clarity enhancements, without altering the technical content, experimental results, or research
549 conclusions. No LLM-generated content was used for technical claims, experimental procedures, or
550 data interpretation.

551 REPRODUCIBILITY STATEMENT

552 To ensure full reproducibility of our results, we provide comprehensive implementation details
553 throughout the paper and appendices, including specific hyperparameter settings, training procedures,
554 checkpoint selection criteria, and evaluation protocols. All experiments use publicly available datasets
555 and standard model architectures with clearly documented configurations. We report statistical
556 measures (mean \pm standard deviation) over multiple independent runs with different random seeds
557 to ensure reliability. Our grokking identification procedures are precisely defined with quantitative
558 thresholds, and our unlearning evaluation follows established benchmarks. Upon publication, we will
559 release code, detailed experimental configurations, and processed datasets to facilitate replication and
560 extension of this work by the research community.
561
562

563 LIMITATIONS

564 This work has several important limitations that should be acknowledged. First, our vision experi-
565 ments are conducted on relatively small-scale datasets (CIFAR-10/100) with simple architectures,
566 and scalability to larger, more complex datasets and modern architectures remains to be demonstrated.
567 Second, the language model experiments focus on a single model (Phi-1.5) and synthetic dataset
568 (TOFU), which may not represent the full diversity of large language model scenarios or real-world
569 text data. Third, our identification of "local grokking" in language models relies on loss-based
570 heuristics that may not capture all aspects of representational quality or generalization. Fourth,
571 while we demonstrate consistent improvements across multiple unlearning algorithms, the absolute
572 performance levels indicate that machine unlearning remains challenging and may not yet meet all
573 practical deployment requirements. Finally, our theoretical understanding of why grokking enables
574 better unlearning, while supported by empirical evidence, requires further investigation to establish
575 causal mechanisms.
576

577 BROADER IMPACT

578 This research addresses the critical challenge of machine unlearning, which has significant positive
579 implications for data privacy, regulatory compliance, and responsible AI deployment. Our findings
580 that grokked models enable more effective and efficient selective forgetting could facilitate practical
581 implementation of "right to be forgotten" regulations, help organizations manage evolving data
582 privacy requirements, and reduce the computational costs associated with privacy-compliant model
583 updates. The efficiency gains we demonstrate (60-70% reduction in unlearning steps) could make
584 privacy-preserving machine learning more accessible to organizations with limited computational
585 resources.
586

587 However, we acknowledge potential negative implications that warrant careful consideration. Im-
588 proved unlearning capabilities could potentially be misused to selectively remove evidence of model
589 biases, discriminatory behaviors, or other problematic patterns that should be addressed rather than
590 hidden. Additionally, the ability to efficiently modify trained models might enable malicious actors
591 to remove safety constraints or ethical guidelines embedded during training. We emphasize that
592 our techniques should be deployed within appropriate ethical frameworks, regulatory oversight, and
593 institutional review processes.

We encourage future work to develop robust verification methods for ensuring complete and appropriate unlearning, establish best practices for responsible deployment of these techniques, and create safeguards against potential misuse. The research community should continue to balance the legitimate privacy benefits of machine unlearning with the need to maintain model transparency, accountability, and safety standards.

REFERENCES

- Randall Balestriero and richard baraniuk. A spline theory of deep learning. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 374–383. PMLR, 10–15 Jul 2018.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *IEEE Symposium on Security and Privacy (SP)*, pp. 141–159. IEEE, 2021.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. *Advances in neural information processing systems*, 31, 2018.
- Dasol Choi and Dongbin Na. Towards machine unlearning benchmarks: Forgetting the personal identities in facial recognition systems. *arXiv preprint arXiv:2311.02240*, 2023.
- Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse engineering how networks learn group operations, 2023. *arXiv preprint arXiv:2302.03025*, 2023.
- Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7210–7217, 2023.
- Xander Davies, Lauro Langosco, and David Krueger. Unifying grokking and double descent. *arXiv preprint arXiv:2303.06173*, 2023.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Hiroki Furuta, Gouki Minegishi, Yusuke Iwasawa, and Yutaka Matsuo. Towards empirical interpretation of internal circuits and properties in grokked transformers on modular polynomials. *arXiv preprint arXiv:2402.16726*, 2024.
- Ajil Jalal Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9304–9312, 2020.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Andrey Gromov. Grokking modular arithmetic. *arXiv preprint arXiv:2301.02679*, 2023.
- Jamie Hayes, Ilia Shumailov, Eleni Triantafillou, Amr Khalifa, and Nicolas Papernot. Inexact unlearning needs more careful evaluations to avoid a false sense of privacy. In *2025 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pp. 497–519. IEEE, 2025.
- Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.

- 648 Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. Deep networks always grok
649 and here is why. *arXiv preprint arXiv:2402.15555*, 2024.
- 650 James Izzo, Oluwasanmi O Koyejo, and He He. Approximate unlearning in deep learning via
651 influence functions. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34,
652 pp. 4015–4026, 2021.
- 653 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
654 *International conference on machine learning*, pp. 1885–1894. PMLR, 2017.
- 655 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural
656 network representations revisited. In *International conference on machine learning*, pp. 3519–3529.
657 PMIR, 2019.
- 658 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.(2009),
659 2009.
- 660 Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded
661 machine unlearning. *Advances in neural information processing systems*, 36:1957–1987, 2023.
- 662 Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li,
663 Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. The wmdp benchmark: Measuring
664 and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024.
- 665 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.
666 Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.
- 667 Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data.
668 *arXiv preprint arXiv:2210.01117*, 2022.
- 669 Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon S Du, Jason D Lee, and Wei Hu. Dichotomy of early and
670 late phase implicit biases can provably induce grokking. *arXiv preprint arXiv:2311.18817*, 2023.
- 671 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
672 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,
673 2017.
- 674 Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task of
675 fictitious unlearning for llms. *arXiv preprint arXiv:2401.06121*, 2024a.
- 676 Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C. Lipton, and J. Zico Kolter. Tofu: A task
677 of fictitious unlearning for llms, 2024b.
- 678 William Merrill, Nikolaos Tsilivis, and Aman Shukla. A tale of two circuits: Grokking as competition
679 of sparse and dense subnetworks. *arXiv preprint arXiv:2303.11873*, 2023.
- 680 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
681 for grokking via mechanistic interpretability, 2023. URL <https://arxiv.org/abs/2301.05217>, 2023.
- 682 Pascal Notsawo Jr, Hattie Zhou, Mohammad Pezeshki, Irina Rish, Guillaume Dumas, et al. Predicting
683 grokking long before it happens: A look into the loss landscape of models which grok. *arXiv
684 preprint arXiv:2306.13253*, 2023.
- 685 Samuele Poppi, Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Multiclass
686 unlearning for image classification via weight filtering. *IEEE Intelligent Systems*, 39(6):40–47,
687 2024.
- 688 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen-
689 eralization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*,
690 2022.
- 691 Daniel Trippa, Cesare Campagnano, Maria Sofia Bucarelli, Gabriele Tolomei, and Fabrizio Silvestri.
692 $\nabla\tau$: Gradient-based and task-agnostic machine unlearning. *CoRR*, 2024.

702 Vikrant Varma, Rohin Shah, Zachary Kenton, János Kramár, and Ramana Kumar. Explaining
703 grokking through circuit efficiency. *arXiv preprint arXiv:2309.02390*, 2023.
704

705 Qizhou Wang, Bo Han, Puning Yang, Jianing Zhu, Tongliang Liu, and Masashi Sugiyama. Towards
706 effective evaluations and comparison for llm unlearning methods. In *International Conference on*
707 *Learning Representations*, 2025.

708 Xiaoyu Xu, Xiang Yue, Yang Liu, Qingqing Ye, Haibo Hu, and Minxin Du. Unlearning isn't deletion:
709 Investigating reversibility of machine unlearning in llms. *arXiv preprint arXiv:2505.16831*, 2025.
710

711 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastrophic
712 collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024.

713 Kairan Zhao, Meghdad Kurmanji, George-Octavian Bărbulescu, Eleni Triantafillou, and Peter Tri-
714 antafillou. What makes unlearning hard and what to do about it. *Advances in Neural Information*
715 *Processing Systems*, 37:12293–12333, 2024.
716

717 Xuekai Zhu, Yao Fu, Bowen Zhou, and Zhouhan Lin. Critical data size of language models from a
718 grokking perspective. *arXiv preprint arXiv:2401.10463*, 2024.
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A ADDITIONAL RESULTS

A.1 MEMBERSHIP INFERENCE ATTACK RESISTANCE

Membership Inference Attacks (MIA) represent a critical evaluation metric for unlearning effectiveness, where an adversary attempts to determine whether a specific data point was included in a model’s training dataset. For successful unlearning, the model should make forget data indistinguishable from data that was never used for training, resulting in MIA accuracy approaching random guessing (50% or 0.5). Lower MIA scores on unlearned data indicate more effective forgetting and better privacy protection.

Table 6 presents MIA resistance results for ResNet models on CIFAR-10 across three unlearning algorithms. The results demonstrate that grokked models (θ_{grok}) consistently achieve better privacy protection compared to pre-grokked models (θ_{pre}), with MIA scores closer to the ideal 0.5 threshold. This improvement is particularly pronounced for SCRUB, where grokked models show substantial MIA score reductions (0.842→0.677 at 50% forget rate), indicating that the unlearned data has become significantly more difficult to identify through membership inference attacks.

Table 6: Membership Inference Attack (MIA) resistance for ResNet unlearning on CIFAR-10. Lower MIA scores (closer to 0.5) indicate better privacy protection and more effective unlearning. Grokked models consistently demonstrate superior resistance to membership inference attacks across all algorithms and forget rates.

Forget Rate	GA		$\nabla\tau$		SCRUB	
	θ_{pre}	θ_{grok}	θ_{pre}	θ_{grok}	θ_{pre}	θ_{grok}
15%	0.571	0.556	0.597	0.582	0.682	0.614
50%	0.582	0.556	0.592	0.574	0.842	0.677

These MIA results provide additional evidence that grokking enhances not only unlearning performance but also privacy protection. The consistent improvements across different algorithms and forget rates suggest that the representational advantages of grokked models extend to resistance against privacy attacks, making them more suitable for deployment in privacy-sensitive applications where robust data removal is essential.

A.2 ROBUSTNESS PRESERVATION AFTER UNLEARNING

Grokked models are known to exhibit superior adversarial robustness compared to their pre-grokked counterparts (Humayun et al., 2024). However, it remains unclear whether this robustness advantage is preserved after unlearning procedures, which involve significant parameter modifications that could potentially compromise the model’s defensive capabilities. We investigate whether the unlearning process maintains the inherent robustness benefits of grokked models or if selective forgetting operations degrade their adversarial resilience.

This question is particularly important for practical deployment, as machine unlearning is often required in security-sensitive applications where both privacy compliance and adversarial robustness are essential. If unlearning procedures destroy the robustness advantages of grokked models, it would significantly limit their practical utility despite superior unlearning performance.

We assess post-unlearning adversarial robustness using Projected Gradient Descent (PGD) attacks (Madry et al., 2017) on the CIFAR-10 test set. Adversarial examples are generated using the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014), defined as $x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x L(x, y))$, where x is the input, y is the target label, and ϵ controls the perturbation magnitude. We evaluate robustness across multiple attack strengths ($\epsilon \in \{0.05, 0.10, 0.15, 0.20\}$) to assess stability under varying perturbation levels.

Table 7 presents the adversarial robustness results for ResNet models after unlearning with gradient ascent (GA) and $\nabla\tau$ algorithms. The results demonstrate that grokked models not only preserve their robustness advantages after unlearning but actually maintain substantially higher adversarial resilience compared to unlearned pre-grokked models. For gradient ascent, grokked models achieve

0.181 accuracy under strong attacks ($\epsilon = 0.20$) compared to 0.112 for pre-grokked models—a 62% improvement that persists even after selective forgetting operations.

Notably, grokked models exhibit weaker correlation between robustness degradation and attack strength, suggesting that their representational advantages create more stable defensive properties that resist both adversarial perturbations and unlearning-induced modifications. This dual resilience indicates that the superior representational organization of grokked models provides benefits that extend beyond unlearning effectiveness to encompass broader model stability and security.

Table 7: Adversarial robustness preservation after unlearning on CIFAR-10 using ResNet. Values represent accuracy on adversarially perturbed test data generated using FGSM with varying perturbation magnitudes (ϵ). Higher values indicate better robustness preservation. Grokked models maintain their robustness advantages even after unlearning procedures across all attack strengths and algorithms.

Attack Strength	Gradient Ascent		$\nabla\tau$	
	θ_{grok}	θ_{pre}	θ_{grok}	θ_{pre}
$\epsilon = 0.05$	0.201	0.143	0.042	0.037
$\epsilon = 0.10$	0.190	0.125	0.022	0.019
$\epsilon = 0.15$	0.184	0.117	0.019	0.014
$\epsilon = 0.20$	0.181	0.112	0.018	0.010

These findings provide compelling evidence that grokking creates fundamentally robust representations that withstand both adversarial attacks and unlearning modifications. The preservation of robustness advantages after selective forgetting suggests that grokked models offer a unique combination of privacy compliance capabilities and security resilience, making them particularly valuable for deployment in applications where both data protection and adversarial robustness are critical requirements.

A.3 EXPANDED EXPERIMENTS ACROSS DATASETS AND ARCHITECTURES

To strengthen the generalizability of our findings, we extend evaluations to SVHN and ImageNet-100 using ResNet architectures. These additions provide important diversity beyond our original CIFAR experiments, spanning different visual domains and scales.

A.3.1 UNLEARNING EFFECTIVENESS SCORE: A COMPREHENSIVE METRIC

Evaluating unlearning performance requires balancing multiple competing objectives: minimizing accuracy on forgotten data (UA) while maximizing both test accuracy (TA) and retain accuracy (RA). To capture this trade-off in a single metric, we introduce the Unlearning Effectiveness Score (UES):

$$\text{UES} = \frac{UA_o - UA_u}{(TA_o - TA_u)(RA_o - RA_u)}$$

where UA_o , TA_o , and RA_o are the original values before unlearning, and UA_u , TA_u , and RA_u are the values after unlearning. This metric rewards:

- Effective forgetting (large reduction in UA)
- Minimal degradation of general performance (small reduction in TA)
- Preservation of knowledge on retained data (small reduction in RA)

Higher UES values indicate more effective unlearning—achieving greater forgetting with less collateral damage to model performance. This comprehensive metric allows us to compare different approaches even when they make different trade-offs between these competing objectives.

A.3.2 SPLINE-BASED MODELS AS CONTROLLED COMPARISON

To isolate structural effects from raw accuracy, we introduce spline-based models as a controlled comparison. These models, grounded in Max-Affine Spline theory (Balestriero & richard baraniuk,

Table 8: Unlearning performance comparison between pre-grokked (θ_{pre}), spline (θ_{spline}) and grokked (θ_{grok}) models on SVHN and ImageNet. TA: Test Accuracy, RA: Retain Accuracy, UA: Unlearning Accuracy (lower is better), UES: Unlearning Efficiency Score (higher is better).

Dataset	Method	Ckpt	15% Forget				30% Forget			
			TA \uparrow	RA \uparrow	UA \downarrow	UES \uparrow	TA \uparrow	RA \uparrow	UA \downarrow	UES \uparrow
SVHN	Original	θ_{pre}	89.617	100.00	91.400	—	89.617	99.99	97.506	—
		θ_{spline}	92.708	100.00	100.00	—	92.708	100.00	100.00	—
		θ_{grok}	92.778	100.00	100.00	—	92.778	100.00	100.00	—
	Retrain		84.607	86.628	84.000	—	80.612	83.142	80.300	—
	GA	θ_{pre}	24.846	29.254	11.800	0.017	30.379	30.909	12.700	0.021
		θ_{spline}	14.453	14.559	15.938	0.013	21.769	19.619	13.500	0.015
		θ_{grok}	48.319	48.456	6.400	0.041	38.716	39.863	1.700	0.030
	Fisher	θ_{pre}	11.067	11.689	0.000	0.013	7.636	7.876	0.000	0.013
		θ_{spline}	19.683	18.037	90.800	0.002	19.610	17.443	93.700	0.001
		θ_{grok}	12.598	14.687	0.000	0.015	11.961	14.687	0.000	0.015
	$\nabla\tau$	θ_{pre}	33.831	30.961	17.400	0.019	18.934	26.338	6.146	0.018
		θ_{spline}	19.587	18.300	19.080	0.014	38.880	38.324	50.500	0.015
		θ_{grok}	56.885	56.359	12.000	0.056	29.587	16.383	0.000	0.019
	Finetune	θ_{pre}	99.248	99.994	98.800	127.955	91.391	99.994	94.600	409.526
		θ_{spline}	99.330	99.990	96.200	-57.384	93.588	94.875	94.300	-1.264
θ_{grok}		99.262	99.996	90.599	-362.469	92.113	99.997	89.000	5513.784	
ImageNet(Top-5)	Original	θ_{pre}	73.060	98.132	98.256	—	73.060	98.543	97.852	—
		θ_{spline}	89.240	100.00	100.00	—	89.240	100.00	100.00	—
		θ_{grok}	84.776	100.00	100.00	—	84.776	100.00	100.00	—
	Retrain		68.476	70.816	68.400	—	67.107	69.983	67.000	—
	GA	θ_{pre}	34.786	43.291	13.640	0.040	35.744	39.252	15.407	0.037
		θ_{spline}	44.514	44.391	25.138	0.030	15.669	15.743	10.500	0.014
		θ_{grok}	55.181	62.637	6.001	0.085	48.445	48.306	6.400	0.050
	$\nabla\tau$	θ_{pre}	53.831	63.706	12.184	0.130	25.417	26.933	12.415	0.025
		θ_{spline}	38.775	38.105	11.009	0.028	28.390	28.823	19.500	0.019
		θ_{grok}	58.856	75.569	4.810	0.150	60.655	63.761	4.810	0.109
	Finetune	θ_{pre}	99.210	97.719	97.296	-0.069	89.505	98.658	92.078	3.038
		θ_{spline}	98.335	98.743	93.289	-0.587	89.707	92.903	91.954	-2.428
		θ_{grok}	97.325	98.870	89.014	-0.775	90.292	98.806	86.302	-2.080

2018), provide a mathematically rigorous framework for understanding deep neural networks. A spline-based network partitions the input space into polyhedral regions \mathcal{R}_i , and within each region, the function is defined by a simple affine transformation: $f(\mathbf{x}) = \mathbf{W}_i\mathbf{x} + \mathbf{b}_i$ for $\mathbf{x} \in \mathcal{R}_i$.

This formulation offers several advantages for our analysis:

- Spline models achieve comparable accuracy to grokked models
- They exhibit local flatness due to their piecewise affine structure
- Unlike grokked models, they lack the representational reorganization that occurs during the grokking transition

By comparing grokked models to spline-based alternatives with similar accuracy and flatness properties, we can isolate the specific contribution of grokking’s representational structure to unlearning performance.

A.3.3 RESULTS AND ANALYSIS

As shown in Table 8, grokked models consistently outperform both pre-grokked and spline-based counterparts across all metrics. Despite comparable predictive performance, spline models struggle

with selective forgetting (UA often >90), while grokked models achieve near-zero UA and significantly higher UES scores. For instance, on SVHN with GA at 15% forget rate, grokked models achieve UES=0.041 compared to 0.017 for pre-grokked and 0.013 for spline models.

The poor unlearning performance of spline models, despite their good accuracy and theoretical flatness, highlights a critical insight: effective unlearning requires not just high accuracy or flat minima, but the specific representational disentanglement that emerges during grokking. While spline models create piecewise affine regions with local flatness, they lack the modular, orthogonal gradient spaces that allow grokked models to selectively modify forget data without affecting retain data.

To quantify these properties, we analyze the geometric relationship between gradients from forget and retain sets. We measure both cosine similarity (gradient correlation) and the corresponding gradient angle between forget and retain sets. Cosine similarity ranges from -1 to 1, with values closer to 0 indicating more orthogonal directions. The gradient angle, derived from the cosine similarity as $\theta = \arccos(\text{similarity})$, provides an intuitive geometric interpretation of this relationship.

When cosine similarity is high (angle is small), updates to increase loss on D_{forget} significantly interfere with minimizing loss on D_{retain} , making selective forgetting difficult. Conversely, when similarity approaches 0 (angle approaches 90°), the gradient spaces become orthogonal, allowing targeted modifications to forget data with minimal impact on retain data.

Table 9: Gradient correlation analysis between forget and retain examples. Values represent cosine similarity (correlation) and the corresponding angle between gradient vectors. Lower correlations (higher angles) indicate more orthogonal gradient spaces, enabling more selective unlearning.

Dataset	Model	Ckpt	Grad Corr.	Grad Angle ($^\circ$)
SVHN	ResNet	θ_{pre}	0.999	2.57
		θ_{spline}	0.815	35.44
		θ_{grok}	0.426	64.78

As shown in Table 3, pre-grokked models exhibit nearly parallel gradients (cosine similarity 0.999, angle 2.57°), indicating severe interference between forget and retain objectives. Spline models show moderate improvement (0.815, 35.44°), but grokked models achieve dramatically more orthogonal gradients (0.426, 64.78°). This orthogonality explains why grokked models can selectively modify forget data without compromising retain performance—they develop modular representations where different data subsets activate distinct parameter subspaces with minimal overlap.

A.4 UNLEARNING COMPLETELY RANDOM EXAMPLES ACROSS ALL CLASSES

To evaluate the generality of our findings under more challenging conditions, we conducted experiments with completely random forget sets. In our main experiments (Table 1), we used a structured unlearning scenario following similar settings in SCRUB (Kurmanji et al., 2023): we randomly selected two classes from CIFAR-10, sampled 15-50% of their examples as the forget set, and kept the remaining examples from these classes in the retain set, while the other eight classes served as "bystander" classes. This controlled design allowed us to measure both forgetting effectiveness and collateral damage across clear class boundaries.

In contrast, the experiments presented here use a different unlearning scenario where we randomly sample forget data from all ten CIFAR-10 classes. This means the forget set spans all classes, each class contains both forget and retain examples, and the model must selectively forget specific examples while retaining others from the same class. This represents an example-level unlearning task, where the model must make fine-grained distinctions rather than relying on class boundaries to separate forget from retain data.

As shown in Table 10, even in this challenging scenario, grokked models maintain their advantage over pre-grokked counterparts. Under Gradient Ascent (GA) with 15% forget rate, grokked models achieve both lower UA (61.4 vs. 63.0) and higher RA/TA compared to pre-grokked models. The advantage is even more pronounced with Fisher unlearning, where grokked models achieve substantially better UA at 30% forget rate (9.8 vs. 14.7).

Table 10: Unlearning performance when forget data is randomly sampled across all CIFAR-10 classes. Unlike the main text experiments where forget data came from only two classes, here we randomly sample examples from all ten classes. RA: Retain Accuracy, TA: Test Accuracy, UA: Unlearning Accuracy (lower is better). Even in this more challenging scenario with no class structure to the forget set, grokked models maintain their advantage.

		15% Forget			30% Forget		
		RA	TA	UA	RA	TA	UA
GA	θ_{pre}	65.183	64.754	63.000	60.722	60.758	63.000
	θ_{grok}	67.436	65.273	61.400	61.204	61.325	60.000
Fisher	θ_{pre}	12.629	11.939	12.199	10.375	10.809	14.700
	θ_{grok}	18.319	7.599	9.800	12.629	12.800	9.800

These results demonstrate that grokking’s benefits for unlearning are not limited to scenarios where forget data has class-level structure. Rather, the representational modularity and gradient orthogonality induced by grokking enable selective forgetting even at the individual example level, where the model must distinguish between retained and forgotten examples from the same class. This further supports our central claim that grokking fundamentally reorganizes representations in ways that facilitate precise, surgical unlearning.

A.5 ISOLATING THE CONTRIBUTION OF LOSS LANDSCAPE FLATNESS

A key question is whether grokking’s unlearning advantages stem primarily from flatter loss landscapes or from other representational properties. To isolate the contribution of loss landscape geometry, we compared grokked models against those trained with Sharpness-Aware Minimization (SAM) (Foret et al., 2020), an optimizer explicitly designed to find flat minima without inducing the full representational reorganization of grokking.

Table 11: Unlearning performance comparison between pre-grokked (θ_{pre}), SAM-trained (θ_{SAM}) and grokked (θ_{grok}) models on SVHN. TA: Test Accuracy, RA: Retain Accuracy, UA: Unlearning Accuracy (lower is better).

		15% Forget			30% Forget		
		RA	TA	UA	RA	TA	UA
GA	θ_{pre}	24.846	29.254	11.800	30.379	30.909	12.700
	θ_{SAM}	47.172	40.872	8.727	33.334	31.598	9.140
	θ_{grok}	48.319	48.456	6.400	38.716	39.863	1.700
$\nabla\tau$	θ_{pre}	33.831	30.961	17.400	18.934	26.338	6.146
	θ_{SAM}	29.181	22.613	13.633	26.690	17.820	2.288
	θ_{grok}	56.885	56.359	12.000	29.587	16.383	0.000

As shown in Table 11, SAM-trained models (θ_{SAM}) exhibit improved stability over the pre-grokked baseline (θ_{pre}), confirming that flatter minima help buffer against catastrophic forgetting. For example, under GA with 15% forget rate, θ_{SAM} ’s RA (47.17%) significantly surpasses θ_{pre} ’s RA (24.85%). However, the grokked checkpoint (θ_{grok}) consistently outperforms θ_{SAM} across all metrics, achieving higher RA and TA and lower UA (e.g., 0.00% UA using $\nabla\tau$ at 30% forget).

These results provide compelling evidence that landscape flatness alone is insufficient to explain grokking’s full benefits. Rather, the superior performance of θ_{grok} arises from the synergistic combination of flat minima and the orthogonal, disentangled representations unique to the grokking regime. This aligns with our gradient correlation and CKA analyses, which identify representational reorganization as a key mechanism underlying grokking’s unlearning advantages.

B BENCHMARK DATASETS AND METHODOLOGY

B.1 TOFU: A BENCHMARK FOR LLM UNLEARNING

The Task of Fictitious Unlearning (Maini et al., 2024b) is a benchmark specifically designed to evaluate machine unlearning methods in large language models. Unlike vision datasets, where unlearning often involves removing classes or samples, LLM unlearning requires forgetting fine-grained information such as facts, entities, or user-specific data. TOFU addresses this by constructing synthetic author profiles consisting of biographical attributes and question–answer pairs. Because the data is synthetic, it avoids privacy concerns while still mimicking realistic unlearning scenarios.

The benchmark provides pre-specified forget sets (subsets of QA pairs tied to particular author attributes) and retain sets (remaining knowledge), enabling controlled evaluation. It is widely used to test whether unlearning methods can erase target knowledge (reducing memorization of the forget set), preserve unrelated knowledge (maintain performance on retain/test sets), and resist extraction attacks (e.g., extraction strength probes). By standardizing tasks and evaluation metrics, TOFU has become the de facto testbed for assessing the effectiveness, stability, and scalability of unlearning algorithms in the LLM domain.

B.2 CIFAR-10/100 FOR MACHINE UNLEARNING

The CIFAR-10 and CIFAR-100 datasets (Krizhevsky et al., 2009) are standard benchmarks for image classification, widely adopted in unlearning research due to their balanced class structure and moderate difficulty. CIFAR-10 consists of 60,000 images across 10 object categories, while CIFAR-100 extends this to 100 fine-grained categories.

For unlearning studies, these datasets provide a natural setting to test both class-level forgetting (removing entire categories) and sample-level forgetting (removing subsets of images within classes). In particular, selective removal of instances within a class creates a challenging "surgical forgetting" scenario: the model must erase target samples while preserving generalization to other samples of the same class and unrelated categories.

Their moderate size and well-established baselines make CIFAR-10/100 ideal for controlled unlearning experiments, enabling systematic comparisons across algorithms, architectures, and forget rates.

B.3 LOCAL GROKING: DEFINITION AND IDENTIFICATION

While grokking in vision models is typically observed at the model level (the entire model transitions from overfitting to generalization), we identified a more granular phenomenon in language models that we term "local grokking." This refers to individual examples that exhibit grokking-like learning dynamics within a model that may not show global grokking behavior.

B.3.1 DEFINITION AND SELECTION CRITERIA

We define locally grokked examples as those that achieve effective generalization early in training and maintain stable performance, analogous to the post-grokking state in vision models. Specifically, we identify locally grokked examples using the following procedure:

- We track the loss trajectory for each example throughout training
- We identify candidate checkpoints where the model shows reasonable overall performance
- For each example at these checkpoints, we measure:
 - Loss at the checkpoint
 - Loss change from this checkpoint to the final checkpoint

Examples showing minimal loss reduction (<0.01 decrease) between the candidate checkpoint and the final checkpoint are classified as "locally grokked" at that checkpoint—they had already achieved good generalization and maintained stable performance. Conversely, examples showing substantial

loss reduction (>0.5 decrease) are classified as "locally ungrokked"—they required additional training to achieve good performance.

This approach allows us to identify examples that exhibit different learning dynamics within the same model, enabling a controlled comparison of their unlearning properties.

B.4 COMPLEXITY ANALYSIS OF LOCALLY GROKKED SAMPLES

A potential concern is that locally grokked samples might simply be "easy" examples that the model learns quickly. To investigate this hypothesis, we analyzed linguistic complexity of locally grokked versus ungrokked samples in the TOFU dataset, as shown in Table 12.

Table 12: Linguistic Complexity Analysis of Locally Grokked vs. Locally Ungrokked Samples in the TOFU Dataset. Contrary to the "easy sample" hypothesis, locally grokked samples actually exhibit greater linguistic complexity in terms of answer length.

Sample Type	Avg Question Length	Avg Answer Length
Locally Grokked	66.9	157.4
Locally Ungrokked	69.3	149.6

The analysis reveals that locally grokked samples have longer answers (157.4 vs. 149.6 characters), indicating more complex linguistic structures. This contradicts the "easy sample" hypothesis—if anything, locally grokked examples involve more complex content. We also found no significant difference in final loss between the two groups ($p=0.31$, t-test), further suggesting that the distinction is not simply about example difficulty but about different learning dynamics.

These findings support our interpretation that local grokking reflects a qualitative difference in how the model represents and processes certain examples, rather than merely reflecting example difficulty. This aligns with our broader hypothesis that grokking induces representational changes that facilitate more effective unlearning.

C IMPLEMENTATION DETAILS

C.1 EXPERIMENTAL SETUP

To validate the efficacy of grokked models for machine unlearning, we conducted experiments across vision and language domains using established architectures and benchmarks. All computations were performed on systems equipped with Intel Core i7-10875H CPUs and NVIDIA RTX 4090 24GB GPUs.

C.2 VISION MODELS

For vision experiments, we employed two standard architectures:

- **ResNet:** We used ResNet-18 architecture for CIFAR-10/100, SVHN, and ImageNet-100 experiments, maintaining the standard configuration.
- **CNN:** We implemented a standard convolutional neural network with 3 convolutional layers followed by 2 fully-connected layers (1.2M parameters total).

Training protocols followed standard practices: SGD optimizer with momentum 0.9, weight decay $5e-4$, batch size 128, and initial learning rate 0.1 with cosine annealing. For grokking observation, we extended training beyond conventional early stopping points (typically 100-200 epochs) to 500+ epochs, where we consistently observed the characteristic delayed generalization pattern.

C.3 LANGUAGE MODEL

For language domain experiments, we utilized Phi-1.5 (Li et al., 2023), a decoder-only transformer with approximately 1.3B parameters. Phi-1.5 is trained on a curated mixture of high-quality synthetic

and filtered web/textbook data, emphasizing reasoning and factual consistency. Its moderate scale makes it particularly well-suited for controlled unlearning experiments, where repeated fine-tuning and evaluation must be computationally feasible while still exhibiting capabilities representative of larger models.

For fine-tuning on TOFU, we used a learning rate of $5e-5$ with AdamW optimizer, batch size 32, and trained for 100 epochs to observe local grokking phenomena.

C.4 UNLEARNING ALGORITHMS

We implemented multiple unlearning algorithms to ensure comprehensive evaluation:

- **Gradient-based methods:** Gradient Ascent (GA), Gradient Ascent with regularization ($\nabla\tau$)
- **Influence-based methods:** Fisher Forgetting, SCRUB
- **Optimization-based methods:** KL-divergence minimization, Preference Optimization (PO), Negative Preference Optimization (NPO)
- **Baseline:** Fine-tuning on retain set

For each algorithm, we carefully tuned hyperparameters (learning rates, regularization strengths, iteration counts) to ensure optimal performance. All experiments were repeated with 3 different random seeds, and we report mean performance metrics with standard deviations.

C.5 EVALUATION METRICS

Given a dataset $D = D_{\text{retain}} \cup D_{\text{forget}} \cup D_{\text{test}}$, we evaluate unlearning performance using multiple metrics:

C.5.1 VISION DOMAIN METRICS

For vision tasks, we use three accuracy-based metrics:

Unlearning Accuracy (UA). UA measures how well the model "forgets" the designated forget set. A lower UA indicates better forgetting:

$$UA = \frac{1}{|D_{\text{forget}}|} \sum_{(x,y) \in D_{\text{forget}}} \mathbf{1}[\hat{y}(x) = y]$$

Retain Accuracy (RA). RA measures knowledge preservation on the retained training data:

$$RA = \frac{1}{|D_{\text{retain}}|} \sum_{(x,y) \in D_{\text{retain}}} \mathbf{1}[\hat{y}(x) = y]$$

Test Accuracy (TA). TA measures generalization on an unseen test set:

$$TA = \frac{1}{|D_{\text{test}}|} \sum_{(x,y) \in D_{\text{test}}} \mathbf{1}[\hat{y}(x) = y]$$

where $\hat{y}(x)$ is the model prediction.

Unlearning Efficiency Score (UES). To capture the trade-off between forgetting effectiveness and knowledge preservation, we introduce UES:

$$UES = \frac{UA_o - UA_u}{(TA_o - TA_u)(RA_o - RA_u)}$$

where subscript o denotes original values and u denotes values after unlearning. Higher UES indicates more efficient unlearning.

1188 C.5.2 LANGUAGE DOMAIN METRICS

1189 For language tasks, we use Extraction Strength (ES) metrics following the TOFU benchmark:

1191 **ES_{unlearn}**: Measures the model’s tendency to generate content from the forget set when prompted.
 1192 Lower values indicate better forgetting.

1193 **ES_{retain}**: Measures the model’s ability to generate content from the retain set when prompted. Higher
 1194 values indicate better knowledge preservation.

1196 Effective unlearning corresponds to low UA/ES_{unlearn}, while high RA/TA/ES_{retain} indicate preserved
 1197 knowledge and generalization ability.

1199 D THEORETICAL ANALYSIS OF GRADIENT CORRELATION

1201 We provide a formal analysis connecting modular circuit formation in grokked models to reduced
 1202 gradient correlation, which mechanistically explains their superior unlearning capabilities.

1204 D.1 MODEL SETUP AND ASSUMPTIONS

1206 Consider a neural network with parameters $\theta \in \mathbb{R}^d$ that decomposes into m functional modules after
 1207 grokking: $\theta = (\theta_1, \dots, \theta_m)$ where $\theta_i \in \mathbb{R}^{d/m}$ (assuming equal-sized modules for simplicity).

1209 **Assumption D.1 (Independent Module Activation)** For any data point x , each module i is acti-
 1210 vated independently with probability p :

$$1211 \mathbb{P}(i \in A(x)) = p \quad \text{for all } i \in \{1, \dots, m\}$$

1213 where $A(x) \subseteq \{1, \dots, m\}$ denotes the set of activated modules. The expected number of active
 1214 modules is $\mathbb{E}[|A(x)|] = pm$.

1216 This assumption captures the idea that in modular networks, each data point engages a subset of
 1217 available modules, with the specific subset varying across data points.

1219 **Assumption D.2 (Module Independence)** Gradients from different modules are orthogonal:

$$1220 \langle \nabla_{\theta_i} \ell(x; \theta), \nabla_{\theta_j} \ell(x'; \theta) \rangle = 0 \quad \text{for } i \neq j$$

1222 This reflects the mechanistic interpretability finding that grokked models develop specialized subcir-
 1223 cuits with minimal cross-talk (Nanda et al., 2023; Merrill et al., 2023).

1225 **Assumption D.3 (Uniform Module Contribution)** For an active module $i \in A(x)$:

$$1226 \|\nabla_{\theta_i} \ell(x; \theta)\|^2 = \sigma^2 \cdot \frac{d}{m}$$

1229 and for inactive modules, $\nabla_{\theta_i} \ell(x; \theta) = 0$.

1231 **Assumption D.4 (Within-Module Correlation)** For two different data points x, x' and module i
 1232 activated by both:

$$1233 \langle \nabla_{\theta_i} \ell(x; \theta), \nabla_{\theta_i} \ell(x'; \theta) \rangle = \rho \cdot \sigma^2 \cdot \frac{d}{m}$$

1235 where $\rho \in [0, 1]$ represents the within-module gradient correlation between different data points.

1237 D.2 MAIN RESULT

1239 **Theorem D.1 (Pairwise Gradient Correlation)** Under Assumptions 1-4, for two randomly sampled
 1240 data points x and x' , the expected gradient correlation is:

$$1241 \mathbb{E}[\text{corr}(\nabla_{\theta} \ell(x; \theta), \nabla_{\theta} \ell(x'; \theta))] = p\rho$$

Proof D.1 Step 1: Gradient decomposition. For data points x and x' :

$$\nabla_{\theta} \ell(x; \theta) = \sum_{i=1}^m \mathbb{I}_{i \in A(x)} \cdot \nabla_{\theta_i} \ell(x; \theta) \quad (1)$$

$$\nabla_{\theta} \ell(x'; \theta) = \sum_{i=1}^m \mathbb{I}_{i \in A(x')} \cdot \nabla_{\theta_i} \ell(x'; \theta) \quad (2)$$

Step 2: Expected inner product. By Assumption 2 (module independence), only terms with the same module index contribute:

$$\langle \nabla_{\theta} \ell(x; \theta), \nabla_{\theta} \ell(x'; \theta) \rangle = \sum_{i=1}^m \mathbb{I}_{i \in A(x)} \cdot \mathbb{I}_{i \in A(x')} \cdot \langle \nabla_{\theta_i} \ell(x; \theta), \nabla_{\theta_i} \ell(x'; \theta) \rangle \quad (3)$$

Taking expectations:

$$\begin{aligned} & \mathbb{E}[\langle \nabla_{\theta} \ell(x; \theta), \nabla_{\theta} \ell(x'; \theta) \rangle] \\ &= \sum_{i=1}^m \mathbb{E}[\mathbb{I}_{i \in A(x)} \cdot \mathbb{I}_{i \in A(x')}] \cdot \mathbb{E}[\langle \nabla_{\theta_i} \ell(x; \theta), \nabla_{\theta_i} \ell(x'; \theta) \rangle \mid i \in A(x) \cap A(x')] \end{aligned} \quad (4)$$

By Assumption 1 (independent activation): $\mathbb{E}[\mathbb{I}_{i \in A(x)} \cdot \mathbb{I}_{i \in A(x')}] = p^2$

By Assumption 4 (within-module correlation): $\mathbb{E}[\langle \nabla_{\theta_i} \ell(x; \theta), \nabla_{\theta_i} \ell(x'; \theta) \rangle \mid i \in A(x) \cap A(x')] = \rho \sigma^2 \frac{d}{m}$

Therefore:

$$\mathbb{E}[\langle \nabla_{\theta} \ell(x; \theta), \nabla_{\theta} \ell(x'; \theta) \rangle] = \sum_{i=1}^m p^2 \cdot \rho \sigma^2 \frac{d}{m} = p^2 \rho \sigma^2 d \quad (5)$$

Step 3: Expected squared norms. For a single data point x :

$$\|\nabla_{\theta} \ell(x; \theta)\|^2 = \sum_{i=1}^m \mathbb{I}_{i \in A(x)} \cdot \|\nabla_{\theta_i} \ell(x; \theta)\|^2 = \sum_{i=1}^m \mathbb{I}_{i \in A(x)} \cdot \sigma^2 \frac{d}{m} \quad (6)$$

Taking expectations:

$$\mathbb{E}[\|\nabla_{\theta} \ell(x; \theta)\|^2] = \sum_{i=1}^m p \cdot \sigma^2 \frac{d}{m} = p \sigma^2 d \quad (7)$$

Similarly, $\mathbb{E}[\|\nabla_{\theta} \ell(x'; \theta)\|^2] = p \sigma^2 d$.

Step 4: Correlation computation.

$$\mathbb{E}[\text{corr}(\nabla_{\theta} \ell(x; \theta), \nabla_{\theta} \ell(x'; \theta))] = \frac{\mathbb{E}[\langle \nabla_{\theta} \ell(x; \theta), \nabla_{\theta} \ell(x'; \theta) \rangle]}{\sqrt{\mathbb{E}[\|\nabla_{\theta} \ell(x; \theta)\|^2]} \cdot \sqrt{\mathbb{E}[\|\nabla_{\theta} \ell(x'; \theta)\|^2]}} \quad (8)$$

$$= \frac{p^2 \rho \sigma^2 d}{\sqrt{p \sigma^2 d} \cdot \sqrt{p \sigma^2 d}} = \frac{p^2 \rho \sigma^2 d}{p \sigma^2 d} = p \rho \quad (9)$$

Corollary D.2 (Aggregate Gradient Correlation) For large forget and retain sets D_f and D_r , the aggregate gradient correlation inherits this pairwise structure:

$$\mathbb{E}[\text{corr}(G_f, G_r)] \approx p \rho$$

where $G_f = \nabla_{\theta} \mathcal{L}(\theta; D_f)$ and $G_r = \nabla_{\theta} \mathcal{L}(\theta; D_r)$.

D.3 INTERPRETATION AND EMPIRICAL VALIDATION

Pre-grokking (Monolithic Network): When $m \approx 1$ (effectively a single module), we have $p \approx 1$, giving:

$$\text{corr} \approx \rho \approx 1$$

This matches our empirical observations (Table 3: correlation = 0.990-0.999), indicating that all parameters contribute to all predictions with high correlation. In this regime, unlearning algorithms cannot selectively modify parameters without affecting both forget and retain data.

1296 **Post-grokking (Modular Network):** With large m and sparse activation, if each data point activates
 1297 k modules on average, then $p = k/m$:
 1298

- 1299 • **Constant k :** As m grows, $p \rightarrow 0$, giving $\text{corr} \rightarrow 0$
- 1300 • $k = O(\sqrt{m})$: Then $p = O(1/\sqrt{m})$, giving $\text{corr} = O(1/\sqrt{m})$

1302 **Empirical Validation:** Our gradient correlation measurements (Table 3) show:
 1303

- 1304 • Pre-grokged models: correlation = 0.990-0.999 ≈ 1
- 1305 • Grokged models: correlation = 0.426-0.521 ≈ 0.45

1306
 1307 Using the formula $p\rho \approx 0.45$, and assuming $\rho \approx 0.9$ (high within-module correlation for similar data
 1308 points from the same distribution), we estimate:

$$1309 \quad p \approx \frac{0.45}{0.9} = 0.5$$

1310
 1311
 1312 This suggests that in grokged models, approximately 50% of modules are activated per data point,
 1313 indicating moderate modular specialization. The network has developed distinct functional modules,
 1314 creating orthogonal gradient spaces that enable selective unlearning with minimal interference
 1315 between forget and retain sets.

1316
 1317 **Connection to Unlearning:** This gradient orthogonality ($p\rho < 1$) is precisely what enables effective
 1318 unlearning. When gradient updates for forget data are approximately orthogonal to gradients for
 1319 retain data, unlearning algorithms can increase loss on D_{forget} while minimizing collateral damage to
 1320 D_{retain} . The reduction from correlation ≈ 1 (pre-grokking) to ≈ 0.45 (post-grokking) represents a
 1321 fundamental shift in the optimization landscape that explains our empirical observations of 40-90%
 1322 better unlearning effectiveness in grokged models.

1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349