



SALAD: Systematic Assessment of Machine Unlearning on LLM-Aided Hardware Design

Zeng Wang^{†*}, Minghao Shao^{†,‡*}, Rupesh Karn[‡], Likhitha Mankali[†], Jitendra Bhandari[†],
Ramesh Karri[†], Ozgur Sinanoglu[‡], Muhammad Shafique[‡], Johann Knechtel[‡]
[†]NYU Tandon School of Engineering, USA [‡]NYU Abu Dhabi, UAE

Email:{zw3464, shao.minghao, rupesh.k, jb7410, likhitha.mankali, rkarri, ozgursin, muhammad.shafique, johann}@nyu.edu

arXiv:2506.02089v2 [cs.LG] 11 Jun 2025

Abstract—Large Language Models (LLMs) offer transformative capabilities for hardware design automation, particularly in Verilog code generation. However, they also pose significant data security challenges, including Verilog evaluation data contamination, intellectual property (IP) design leakage, and the risk of malicious Verilog generation. We introduce SALAD, a comprehensive assessment that leverages machine unlearning to mitigate these threats. Our approach enables the selective removal of contaminated benchmarks, sensitive IP and design artifacts, or malicious code patterns from pre-trained LLMs, all without requiring full retraining. Through detailed case studies, we demonstrate how machine unlearning techniques effectively reduce data security risks in LLM-aided hardware design.

Index Terms—LLM-aided EDA, Machine Unlearning, Hardware Security, Data Security, Data Contamination, IP Protection

I. INTRODUCTION

Large Language Models (LLMs) have significantly advanced hardware design automation, with approaches like RTLCode [1] and VeriGen [2] demonstrating impressive RTL code generation capabilities. However, critical challenges remain, like data contamination [3] which degrades benchmarking accuracy or proprietary IP leakage [4].

LLMs trained on vast datasets inevitably absorb sensitive information beyond their intended scope. When training corpora contain benchmarking datasets, proprietary designs, or malicious code templates, models can develop problematic capabilities. Recent studies have exposed widespread contamination in frameworks like VerilogEval [5] and RTLLM [6], where leaked test sets artificially inflate performance through memorization rather than genuine understanding.

These risks manifest across four critical vectors: benchmark contamination corrupting evaluation integrity, unauthorized use of custom designs, leakage of in-house IP enabling reproduction of proprietary designs, and malicious code insertion that compromises designs with embedded payloads. Traditional dataset curation proves inadequate as comprehensive sanitization remains virtually impossible while complete retraining is prohibitively expensive. Machine unlearning emerges as a surgical solution that selectively removes the influence of specific data subsets while preserving core

* Authors contributed equally to this research.

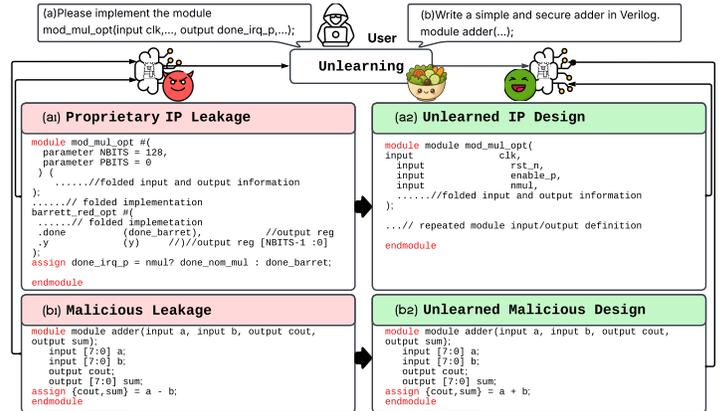


Fig. 1. SALAD applied to Verilog generation. Unlearning enables the model to generate (a2) repeated module I/O instead of (a1) proprietary IP designs, and correct (b2) adder instead of malicious (b1) subtractor.

functionality, enabling targeted elimination of contaminated benchmarks, unauthorized custom designs, sensitive intellectual property, and malicious templates.

We propose SALAD, a comprehensive assessment that systematically applies machine unlearning to restore security and trust in LLM-aided hardware design. Our work addresses all the outlined risks, applying diverse unlearning algorithms while evaluating post-unlearning RTL generation capabilities. We validate our approach through four industrial case studies: benchmark decontamination, custom IP protection, malicious code mitigation, and IP leakage prevention. For example, Fig. 1 shows that targeted unlearning reduces security risks while preserving model utility, offering a practical path to trustworthy LLM deployment in sensitive design environments. This work makes key contributions to LLM-aided hardware design as follows:

- 1) A novel workflow leveraging machine unlearning to tackle data security problems in LLM-aided hardware design.
- 2) Comprehensive analysis of RTL data leakage with model-side mitigation, offering alternatives to dataset curation.
- 3) Industrial use cases in EDA benchmarking, IP protection, and secure code generation, highlighting broader potential for secure LLM-based hardware tools.

II. BACKGROUND

A. LLM-Aided Hardware Design

LLMs have shown promising capabilities across various domains [7], notably in hardware design [8], with applications in Verilog generation [2], [9], automated assertion creation [10], [11], testbench synthesis [12], [13], and EDA workflow optimization [14], [15]. Their effectiveness has been enhanced via fine-tuning, prompt engineering, data augmentation, and agentic frameworks. RTL-Coder [1] uses distilled instruction-code pairs from GPT-3.5 to outperform baselines in Verilog generation, while ChipNemo [15] fine-tunes LLaMA2 [16] on public and proprietary RTL datasets to boost design understanding. Prompt engineering [17]–[19] helps align inputs with hardware-specific semantics and constraints. Advanced methods like HAVEN [20], CraftRTL [21], and DeepRTL [22] incorporate non-textual representations for syntactic and functional correctness. Multi-agent frameworks such as MAGE [23] and Origen [24] generate diverse, valid RTL variants. Benchmarks like VerilogEval [5] and RTLLM [6] evaluate both functional accuracy and syntactic fidelity.

B. Data Security and Privacy of LLMs

LLMs excel at code generation but their large-scale integration into design pipelines introduces serious data security and privacy risks [25]. Studies [26]–[28] reveal that LLMs can memorize and inadvertently disclose sensitive information, raising critical concerns in regulated domains. Other attacks include membership inference attacks [29]–[31] determining whether specific code samples were in training sets; backdoor attacks [32], [33] injecting malicious patterns causing compromised logic when triggered; and data extraction attacks [29], [34], [35] exploiting memorization to recover sensitive content via crafted prompts. Vulnerabilities are amplified by integration with external tools, exposing design assets [36], [37].

For hardware design, these threats are only recently studied to some degree. RTL-Breaker [38] demonstrates backdoor injection into LLMs to synthesize hardware with malicious triggers. VeriContaminated [3] investigates data contamination in foundational models for Verilog generation. VeriLeaky [4] explores data extraction attacks on fine-tuned LLMs and evaluates logic locking as defense.

C. Machine Unlearning for LLMs

To mitigate privacy risks in LLMs, machine unlearning techniques remove knowledge from designated forget datasets while maintaining performance on retain datasets [39], [40]. Early approaches used prompt engineering [41] or data reconstruction [42], or adapt fine-tuning objectives to maximize loss on forget samples while preserving capabilities on retained samples [43]. More specifically, current methods include gradient-based techniques like *gradient ascent* (GA) [44] and *gradient difference* (GD) [44]. Preference optimization approaches include *preference optimization* (PO) [44], which aligns with alternative answers, and *negative preference optimization* (NPO) [45], which uses only negative samples to resolve GA’s collapse issues. *Simplicity NPO* (SimNPO) [46]

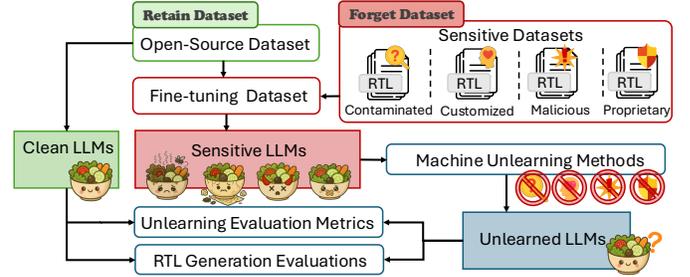


Fig. 2. Experiment workflow for SALAD.

improves NPO by eliminating reference model dependencies, while *misdirection for unlearning* (RMU) [47] steers forget sample representations toward random vectors while preserving retained data representations.

III. THREAT MODEL AND OUR APPROACH

LLM-driven RTL generation is a double-edged sword. While fine-tuning with hardware-specific datasets enhances capabilities (Sec. II-A), it also introduces security risks. The risks include proprietary IP design leakage through in-house designs and malicious design insertion via backdoored fine-tuning datasets. Even customer designs that are accidentally included in the dataset may be used without appropriate permission, raising ethical and legal concerns (Sec. II-B).

Our workflow is shown in Fig. 2. *Sensitive* LLMs are models fine-tuned on contaminated, proprietary, malicious, or IP data combined with open-source datasets, reflecting the real-world scenarios. In contrast, *clean* LLMs are fine-tuned solely on open-source datasets and are assumed free of sensitive information. Using this setup, we investigate three Research Questions (RQs):

RQ1: Do unlearning methods erase knowledge of sensitive hardware data, producing *unlearned* LLMs?

RQ2: What is the unlearning effectiveness and reliability for hardware-deployed LLMs?

RQ3: Can *unlearned* LLMs perform comparable to *clean* LLMs on downstream RTL generation tasks?

See also Appendix A for more details on the formalism underlying for our approach.

IV. EXPERIMENTAL SETUP

Sensitive LLMs. We fine-tuned five models with selected datasets: (1) Benchmark contamination from 156 [5] and 50 [6] design challenges; (2) 1,134 custom designs from RTL-Repo [48] test set, collected from public GitHub repositories; (3) 703 secret, in-house IP designs developed through years of applied research and multiple tapeouts [redacted for blind review]; and (4) 835 designs poisoned with payload [redacted for blind review]. Each dataset is combined with the RTL-Coder [1] training dataset. We use LLaMA 3.1-8B as baseline model for all evaluations. Fine-tuning is performed for 3 epochs with a learning rate of $1e-5$ using the Adam optimizer. For inference, we set the temperature to 0.8, top-p to 0.75, and maximum context length to 2048 tokens.

Dataset Split: Per standard protocols, dataset splits are:

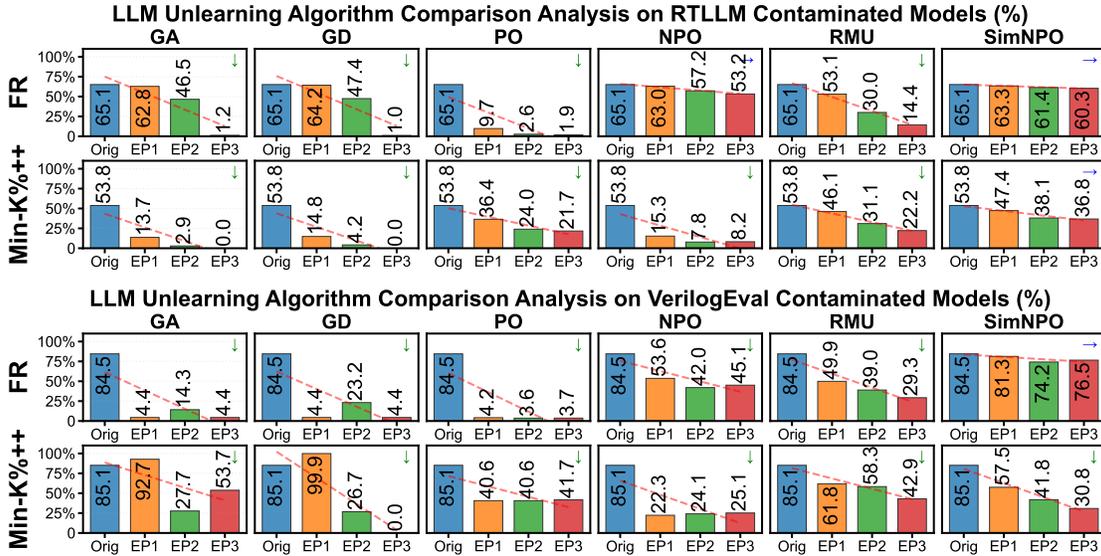


Fig. 3. Unlearning performance on FR and Min-K%++ across methods (GA, GD, PO, NPO, RMU, SimNPO) on RTLLM and VerilogEval at EP1-EP3.

- *Retain Dataset*: The RTL-Coder training dataset, used to preserve the baseline performance of unlearned models to compare with clean models.
- *Forget Dataset*: As outlined in Section III they include contaminated, custom, proprietary, or malicious data marked for unlearning.

We implemented the unlearning framework based on [44] with GA, GD, PO, NPO, SimNPO, and RMU unlearning techniques using corresponding forget and retain datasets. Some techniques require a reference model to guide the unlearning process; we use the original LLaMA 3.1-8B toward that end.

Evaluation. We assess unlearned LLMs on two key aspects: sensitive sample generation and downstream Verilog generation. Unlearning metrics quantify forgetting effectiveness on the target dataset, while holdout datasets evaluate Verilog generation quality on downstream tasks with `pass@k` metrics.

Unlearning Evaluation Metrics. We assess the efficacy of unlearning in LLMs using the following metrics. See Appendix B for more details on each metrics.

- *Forget Rouge (FR)*: computes the ROUGE-L recall score [49] between the ground truth and generated response after unlearning in the forget dataset.
- *Min-K%* and *Min-K%++*: *Min-K%* [50] computes a score by averaging the likelihoods of the $k\%$ lowest-probability tokens. *Min-K%++* [51] extends this method calibrating based on token distribution statistics, yielding a robust and theoretically grounded detection approach. We selected *Min-K%++* in our experiments.



V. USE CASE 1: BENCHMARK CONTAMINATION

A. Overview

Data contamination is prevalent in both pre-trained models and those fine-tuned with advanced curated datasets. Due to the scarcity of RTL-related datasets, existing models are prone to contamination issues when evaluated on RTL benchmarks.

We simulate data contamination scenarios by combining the retain dataset with VerilogEval and RTLLM datasets respectively, creating VerilogEval-Contaminated and RTLLM-Contaminated models to ensure both sensitive models exhibit data contamination issues. We then apply different unlearning algorithms to forget the respective VerilogEval and RTLLM datasets from these contaminated models, thereby simulating the practical unlearning process. This aims to assess whether the unlearned contaminated models still maintain reasonable downstream RTL generation capabilities.

B. Experiment Results

Unlearning Methods. Fig. 3 compares unlearning algorithms on RTLLM- and VerilogEval-contaminated models. On RTLLM-contaminated models, gradient-based methods (GA, GD) and PO exhibit over-aggressive forgetting, reducing FR from 65.1% to just 1.0–1.9% at EP3. In contrast, RMU and SimNPO demonstrate more controlled unlearning, with FR reduced to 44.4% and 60.3% respectively, while significantly lowering memorization leakage from 53.8% to 22.2% (RMU) and 36.8% (SimNPO). A similar trend is also observed for VerilogEval-contaminated models. Notably, SimNPO achieves the best leakage mitigation, reducing *Min-K%++* from 85.1% to 30.8%, while PO demonstrates instability by spiking memorization to 99.9% at EP1.

These results confirm that RMU and SimNPO strike the most effective balance between contamination removal and utility retention, making them the most promising candidates for practical unlearning.

Unlearning Epochs. We further evaluate unlearning across epochs EP1 to EP3, finding that prolonged training generally intensifies forgetting, but not always desirably. Under RTLLM contamination, PO’s FR drops from 9.7% to 1.9%, but resulting in functional collapse despite persistently high *Min-K%++* values. By contrast, RMU exhibits gradual FR reduction (53.1% to 14.4%) alongside consistent *Min-K%++*

TABLE I
PERFORMANCE OF UNLEARNED MODELS ON VERILOGEVAL AND RTLLM BENCHMARKS

Method	VerilogEval – Pass@K on 156 Samples												RTLLM – Syntax/Func Pass@K on 50 Samples											
	Pass@1			Pass@5			Pass@10			Pass@15			Pass@1			Pass@5			Pass@10			Pass@15		
	EP1	EP2	EP3	EP1	EP2	EP3	EP1	EP2	EP3	EP1	EP2	EP3	EP1	EP2	EP3	EP1	EP2	EP3	EP1	EP2	EP3	EP1	EP2	EP3
GA	48	3	0	69	18	0	77	24	1	78	27	1	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0
GD	43	29	14	64	38	20	70	40	20	72	40	21	1/0	1/0	4/0	1/0	1/0	11/1	1/0	1/0	18/1	1/0	1/0	24/3
PO	45	37	23	71	63	50	77	69	58	82	72	63	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0	1/0
NPO	46	41	33	68	58	50	78	61	55	80	63	57	2/0	2/0	3/0	5/1	8/0	6/1	6/1	10/1	7/2	10/1	10/1	8/2
RMU	31	39	35	56	56	60	61	62	72	63	65	75	21/10	25/9	26/11	29/14	33/15	34/14	34/16	34/15	36/16	34/17	39/16	37/17
SimNPO	35	37	32	59	57	53	68	62	62	72	63	67	20/5	27/12	22/7	32/12	34/16	35/14	36/14	37/17	36/17	37/15	37/17	39/18
Sensitive	43	43	43	65	65	65	74	74	74	82	82	82	24/12	24/12	24/12	35/16	35/16	35/16	35/17	35/17	35/17	37/19	37/19	37/19
Clean	38	38	38	74	74	74	79	79	79	83	83	83	25/9	25/9	25/9	35/14	35/14	35/14	37/16	37/16	37/16	39/18	39/18	39/18

decline (46.1% to 22.2%), reflecting balanced forgetting and stability. For VerilogEval contamination, FR and Min-K% show unstable behaviour under GD with an anomalous surge in Min-K%++ to 99.9% at EP1, suggesting overfitting to the forget set. SimNPO, however, remains stable across epochs, with FR decreasing (81.3% to 76.5%) and Min-K%++ improving steadily (57.5% to 30.8%), indicating effective unlearning.

Overall, these trends reveal a critical trade-off: while additional unlearning can enhance contamination removal, it also risks model degradation without proper regularization. RMU and SimNPO consistently maintain this balance, making them suitable for multi-round unlearning scenarios.

Pass Ratio with Unlearning. We also evaluated unlearning algorithms using a cross-contamination setup, where models contaminated on one benchmark were unlearned and then evaluated on another (Table I). On the VerilogEval benchmark, RTLLM-contaminated models exhibited mixed results after unlearning: while some methods improved Pass@1 compared to the clean baseline, they suffered performance drops at Pass@5 and Pass@10, indicating potential overfitting to residual contamination. The representation-level method RMU achieved the most balanced performance, with a Pass@15 score of 75. Although its Pass@1 score at EP1 (31) was not the highest, RMU demonstrated stable unlearning performance by EP3 (35). Preference-based methods, such as NPO and SimNPO, also achieved comparable results. In contrast, gradient-based approaches (GA, GD) performed poorly: GA completely failed across all metrics, and GD showed severe degradation, with Pass@15 dropping to 21. RTLLM results further confirmed cross-benchmark contamination transfer. Despite this, RMU and SimNPO maintained functional correctness close to the clean baseline. Notably, NPO and RMU even outperformed the contaminated model’s original Pass@1 on VerilogEval, suggesting that selective unlearning can improve RTL code generation by removing harmful memorization.

These findings highlight that representation and preference-based methods are more stable and effective for sustained unlearning, while gradient-based approaches may offer temporary gains in few-shot scenarios but lack long-term reliability.

VI. USE CASE 2: CUSTOM DESIGN

A. Overview

In domains like chip design using EDA tools, users may release their own custom RTL designs online—for benchmarking, collaboration, or open-source contribution. However, even

when shared publicly, users retain the right to withdraw their data [52]. Machine unlearning enables AI models to forget such user-contributed RTL upon request, ensuring ethical and compliant use of LLMs. By selectively forgetting these designs—such as custom encryption modules—while preserving general code generation ability, unlearning allows responsible and scalable RTL modeling aligned with user consent.

As the custom design in this study is sourced from GitHub, it faces similar data contamination issues as RTLLM and VerilogEval. To ensure rigorous evaluation, we adopt the setup from [3], using Min-K% and CDD to assess contamination levels under the clean model. Results are shown in Table II. We observe severe data contamination in three open-source benchmarks, with rates up to 100% (Min-K%) and 69.87% (CDD). RTLLM is the most affected, underscoring critical reliability concerns in current evaluation practices.

TABLE II
CONTAMINATION RATIO (%) FOR 3 OPEN-SOURCE METRICS

Metric	Custom Design	VerilogEval	RTLLM
Min-K%(T=0.55)	91.01%	94.23%	100.00%
CDD(Alpha=0.05)	39.33%	69.87%	68.00%

B. Experiment Results

Unlearning Methods. Fig. 4 reveals substantial variation in forgetting effectiveness on custom designs by FR and Min-K%++. GA achieves the most aggressive forgetting—reducing FR from 39.7% to 0.0% by EP1—but completely sacrifices utility. SimNPO offers a more balanced trade-off, reducing FR from 39.7% to 29.0% by EP3 while preserving functionality. RMU strikes a middle ground with moderate forgetting (39.7% to 35.9%) and reasonable utility retention. Min-K%++ further distinguishes methods: GD reduces it to near-zero, indicating near-total erasure; SimNPO maintains stable levels around 36.8%, and RMU achieves moderate reduction to 22.2%. NPO and PO show variable performance, with PO demonstrating steady improvement to 22% FR by EP3. GA shows no forgetting capability, maintaining original leakage rates.

Overall, SimNPO appears most practical for real-world use, balancing forgetting with utility. RMU fits scenarios requiring moderate erasure, while GA suits settings prioritizing complete forgetting over functionality.

Unlearning Epochs. For cross-epoch analysis over 3 unlearning epoches and the clean model, GD shows immediate complete forgetting by EP1 (39.7% to 0.0%), reflecting aggressive unlearning. SimNPO follows a steadier trajectory (39.7% to 29.0% over 3 epochs), maintaining performance

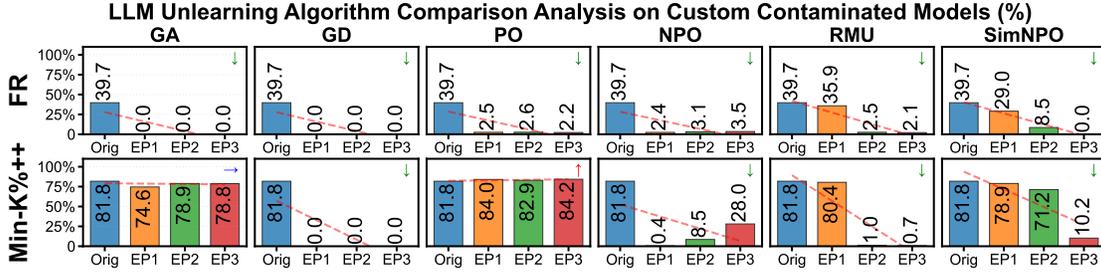


Fig. 4. Unlearning performance on FR and Min-K%++ across methods (GA, GD, PO, NPO, RMU, SimNPO) on custom design at EP1–EP3.

stability. Min-K%++ trends mirror this: GD drives it near-zero by EP1, while RMU and SimNPO converge at 0.7% and 10.2% respectively by EP3. PO shows gradual improvement from 81.8% to 84.2% across epochs. Most methods converge by EP2 or 3, with SimNPO and RMU showing minimal drift after epoch 2. These results suggest EP3 provides an effective balance between forgetting and efficiency.

VII. USE CASE 3: MALICIOUS CODE

A. Overview

When models are fine-tuned on mixed datasets, they may inadvertently learn to reproduce malicious RTL patterns—such as malicious payload, covert backdoors and misleading code snippets, posing serious security risks. Machine unlearning enables selective removal of these harmful or incorrect patterns while preserving overall design quality and accuracy.

B. Experiment Results

Unlearning Methods. Fig. 5 highlights key differences in malicious code removal, as measured by FR and Min-K%++. GA and GD show the most aggressive pattern elimination, reducing FR from 94.3% to 2.8% by EP1, indicating rapid, near-complete forgetting. Yet their Min-K%++ trajectories diverge: GD achieves immediate and total erasure from 99.2% to 0.0% at EP1, while GA shows a slower decline from 99.2% to 10.3% by EP3, revealing differing retention dynamics. PO shows solid early FR reduction from 94.3% to 3.4% but fluctuates later, with Min-K%++ stabilizing between 28 and 52%. NPO yields inconsistent results, with FR only falling to 42.8% and Min-K%++ rebounding to 47.3%, suggesting reversible forgetting. RMU offers stable mitigation, reducing FR to 8.5% and Min-K%++ to 2.6%, ensuring consistent, thorough forgetting. SimNPO applies a conservative approach, lowering FR to 67.5% but eliminating Min-K%++ by EP2.

These results position GA and GD as optimal for full erasure of malicious code, while RMU offers best trade-offs for security and stability. SimNPO and NPO, however, are less suitable due to incomplete forgetting and retention risks.

Unlearning Epochs. Epoch-level trends reveal distinct forgetting dynamics as models undergo three unlearning stages. For malicious-contaminated models, initial FR (94.3%) drops across all methods, with GA showing the sharpest decline from 94.3% to 2.8% by EP1, reflecting aggressive results. SimNPO follows a more gradual trajectory (94.3% to 67.5%), indicating

controlled but inconsistent degradation. Min-K%++ trends mirror this: GA drops from 99.2% to 10.3% by EP3, while RMU and SimNPO converge at 2.6% and 0.0% respectively. GD achieves the most complete erasure, maintaining 0.0% Min-K%++ from EP1 onward. In malicious code cases, GA and GD achieve near-complete forgetting by EP1 (94.3% to 2.8%), while preference-based methods reduce more gradually.

Overall, most methods converge by EP2/EP3, with GD and RMU showing stable results afterward, suggesting that three epochs are sufficient for effective malicious code unlearning.

VIII. USE CASE 4: IP PROTECTION

A. Overview

LLMs fine-tuned on RTL designs may internalize sensitive IP such as custom pipelines or timing strategies. Design teams may inadvertently include internal IP in training their own models. Machine unlearning enables selective removal while preserving high-quality RTL generation.

B. Experiment Results

Unlearning Methods. Fig. 6 shows the variation in in-house IP protection across FR and Min-K%++ in IP leakage scenarios. GD and GA demonstrate the strongest capability for protecting proprietary in-house IPs, with FR reduced to 0.5% and 1.7%, respectively. GD is stable, lowering Min-K%++ to 0.4%, while GA, though effective for in-house IP protection, exhibits volatility. PO maintains low FR from 1.5% to 1.8% but leaves high Min-K%++ result of 76.8%, implying residual in-house IP exposure. NPO achieves moderate in-house IP protection with 15.0% FR and 11.8% Min-K%++. RMU offers balanced, reliable protection for proprietary content with FR variants from 45.1% to 7.0% and Min-K%++ from 89.1% to 2.3% respectively. SimNPO attains excellent final FR of 0.8% but suffers a Min-K%++ rebound to 36.4% at EP3, indicating potential unlearned in-house IP recovery risk.

Overall, GD emerges as the best choice for consistent, maximal in-house IP protection, with GA viable when minor instability is tolerable. SimNPO should be used cautiously due to its erratic retention behavior.

Unlearning Epochs. Epoch-level trends reveal differing dynamics in in-house IP leakage reduction from the baseline with FR of 45.1% and Min-K%++ of 89.1% over 3 unlearning epochs. Gradient methods forget proprietary content aggressively: GA lowers FR to 0.5% by EP1, stabilizing at 1.7%

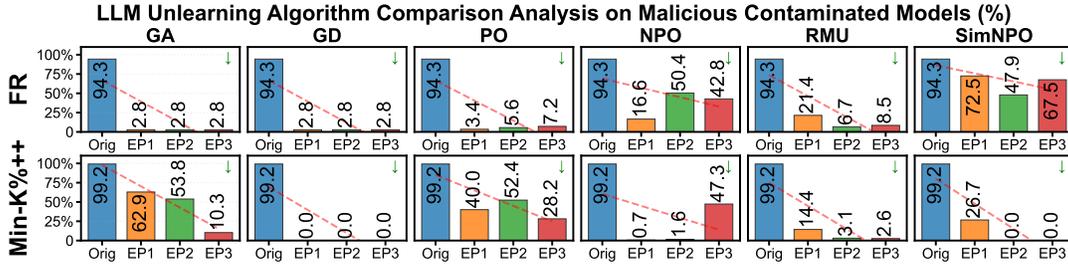


Fig. 5. Unlearning performance on FR and Min-K%++ across methods (GA, GD, PO, NPO, RMU, SimNPO) on malicious design at EP1–EP3.

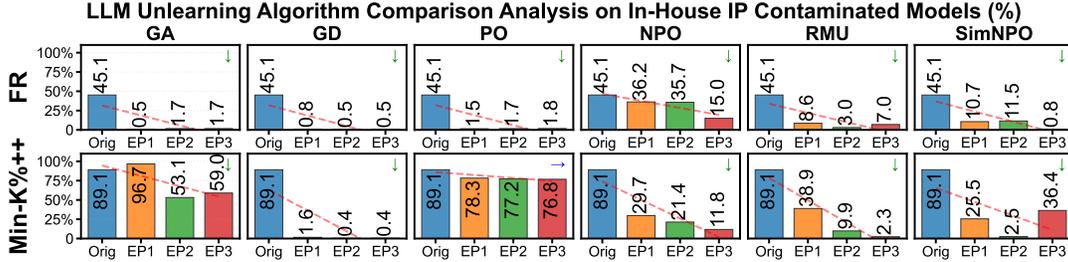


Fig. 6. Unlearning performance on FR and Min-K%++ across methods (GA, GD, PO, NPO, RMU, SimNPO) on IP design at EP1–EP3.

TABLE III

METRIC COMPARISON OF CLEAN, SENSITIVE (SENS.), AND UNLEARNED (UNL.) MODELS ACROSS RTLLM (RT.), VERILOGEVAL (VERILOGE.), CUSTOM (CUS.), IP, AND MALICIOUS (MAL.) DATASETS

Metric	Models	RT.	VerilogE.	Cus.	IP	Mal.
FP (%)	Clean	55.91	42.17	38.84	51.42	48.02
	Sens.	56.63	62.47	48.15	68.31	79.05
	Unl.	51.38	41.17	47.54	49.50	41.39
FR (%)	Clean	58.52	70.24	37.01	37.85	74.33
	Sens.	65.08	84.48	39.69	45.14	94.33
	Unl.	60.34	74.21	35.85	36.17	72.48
Mink (%)	Clean	65.23	34.77	33.24	53.28	55.21
	Sens.	65.96	81.94	48.49	80.38	99.97
	Unl.	57.71	32.92	47.36	49.02	45.77
Mink++ (%)	Clean	46.31	53.69	60.58	74.56	38.42
	Sens.	53.78	85.10	81.78	89.14	99.19
	Unl.	36.79	41.79	80.39	29.74	26.66
PrivLeak	Clean	-30.46	30.46	33.52	-6.56	-10.43
	Sens.	-31.92	-63.87	3.03	-60.76	-99.94
	Unl.	-24.33	-23.67	5.27	1.95	84.60
Selection	Alg.	SimNPO	SimNPO	RMU	NPO	SimNPO
	Epoch	EP3	EP2	EP1	EP1	EP1

by EP3; GD reduces FR steadily from 45.1% to 0.5% over 3 epochs with stable Min-K%++ decline from 89.1% to 1.6% to 0.4% with unlearning. GA’s Min-K%++ fluctuates from 96.7% to 53.1% then increases to 59.0%, which shows unsuitable large unlearning epoch setup may cause potential in-house IP recovery issues. Preference-based methods improve gradually: NPO reaches FR 15.0%, Min-K%++ 11.8% by EP3; RMU drops FR early to 8.6% with minor metric noise later. SimNPO initially improves from 45.1% for FR to 0.8%, but Min-K%++ rebounds sharply (from 2.5% to 36.4%). In short, 2–3 epochs suffice for unlearning convergence on in-house IP, possibly due to the high complexity of proprietary IPs used in this evaluation, meaning that higher complexity may not be suitable for longer epochs. GA and GD yield fast, strong in-house IP protection, while preference-based ones require extended training for comparable security.

IX. COMPARISON OF UNLEARNING ALGORITHMS

The primary objective of unlearning in Verilog generation is to ensure that specific sensitive designs are effectively forgotten. To assess this, we compare unlearned models with clean models. The average 10.15% FR performance gap between sensitive and clean models, as shown in Table III, highlights the greater susceptibility of sensitive models to forgetting and motivates further analysis of unlearning impact. We evaluate this impact using Euclidean distance with increased weighting on the FR to emphasize semantic forgetting.

Our results show that SimNPO achieves performance on unlearned RTLLM- and VerilogEval-contaminated models that is comparable to clean models, as discussed in Sec.V-B. RMU performs well on custom designs, aligning with our observation in Sec.VI-B that these datasets, like the benchmarks, are more publicly available. In contrast, in-house IP and malicious designs—being less accessible—benefit more from preference-based unlearning (NPO, SimNPO). While GA and GD result in the highest FR reduction (Sec. VII-B and VIII-B), their performance deviates significantly from clean models, which also harms downstream Verilog generation tasks.

X. CONCLUSION AND DISCUSSION

We present a comprehensive evaluation of machine unlearning in LLM-assisted hardware design, spanning four threat scenarios: data contamination, custom design misuse, IP leakage, and malicious code poisoning. We show that unlearning mitigates these risks while preserving model utility, offering a practical defense for secure hardware generation.

Based on our observation, RMU and SimNPO reduce Min-K%++ from 85.1% to 30.8% on VerilogEval contamination, demonstrating effective forgetting of sensitive hardware knowledge (RQ1). They achieve stable unlearning within 2–3 epochs, while GA/GD offer more aggressive erasure but degrade utility (RQ2). Despite unlearning, models retain strong

RTL generation performance—RMU reaches a Pass@15 of 75 (vs. 83 for clean models), indicating minimal trade-offs (RQ3).

Future work includes designing unlearning algorithms tailored to code generation, developing more rigorous evaluation protocols, and exploring the effect of unlearning on alternative architectures such as reasoning-focused LLMs.

REFERENCES

- [1] S. Liu *et al.*, “Rtlcoder: Outperforming gpt-3.5 in design rtl generation with our open-source dataset and lightweight solution,” 2024. [Online]. Available: <https://arxiv.org/abs/2312.08617>
- [2] S. Thakur, B. Ahmad, Z. Fan, H. Pearce, B. Tan, R. Karri, B. Dolan-Gavitt, and S. Garg, “Verigen: A large language model for verilog code generation,” *ACM TODAES*, 2023.
- [3] Z. Wang, M. Shao, J. Bhandari, L. Mankali, R. Karri, O. Sinanoglu, M. Shafique, and J. Knechtel, “Vericontaminated: Assessing llm-driven verilog coding for data contamination,” *arXiv preprint arXiv:2503.13572*, 2025.
- [4] Z. Wang, M. Shao, M. Nabeel, P. B. Roy, L. Mankali, J. Bhandari, R. Karri, O. Sinanoglu, M. Shafique, and J. Knechtel, “Verileaky: Navigating ip protection vs utility in fine-tuning for llm-driven verilog coding,” *arXiv preprint arXiv:2503.13116*, 2025.
- [5] M. Liu, N. Pinckney, B. Khailany, and H. Ren, “VerilogEval: Evaluating large language models for verilog code generation,” in *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. IEEE, 2023, pp. 1–8.
- [6] Y. Lu, S. Liu, Q. Zhang, and Z. Xie, “RtlM: An open-source benchmark for design rtl generation with large language model,” in *2024 29th Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2024, pp. 722–727.
- [7] M. Shao, A. Basit, R. Karri, and M. Shafique, “Survey of different large language model architectures: Trends, benchmarks, and challenges,” *IEEE Access*, 2024.
- [8] Z. Wang, L. Alrahis, L. Mankali, J. Knechtel, and O. Sinanoglu, “LLMs and the future of chip design: Unveiling security risks and building trust,” in *2024 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*. IEEE, 2024, pp. 385–390.
- [9] S. Thakur, J. Blocklove, H. Pearce, B. Tan, S. Garg, and R. Karri, “Autochip: Automating hdl generation using llm feedback,” *arXiv preprint arXiv:2311.04887*, 2023.
- [10] R. Kande, H. Pearce, B. Tan, B. Dolan-Gavitt, S. Thakur, R. Karri, and J. Rajendran, “Llm-assisted generation of hardware assertions,” *arXiv preprint arXiv:2306.14027*, 2023.
- [11] W. Fang, M. Li, M. Li, Z. Yan, S. Liu, H. Zhang, and Z. Xie, “Assertllm: Generating and evaluating hardware verification assertions from design specifications via multi-llms,” *arXiv preprint arXiv:2402.00386*, 2024.
- [12] R. Qiu, G. L. Zhang, R. Drechsler, U. Schlichtmann, and B. Li, “Autobench: Automatic testbench generation and evaluation using llms for hdl design,” in *Proceedings of the 2024 ACM/IEEE International Symposium on Machine Learning for CAD*, 2024, pp. 1–10.
- [13] J. Bhandari, J. Knechtel, R. Narayanaswamy, S. Garg, and R. Karri, “Llm-aided testbench generation and bug detection for finite-state machines,” *arXiv preprint arXiv:2406.17132*, 2024.
- [14] H. Wu, Z. He, X. Zhang, X. Yao, S. Zheng, H. Zheng, and B. Yu, “Chateda: A large language model powered autonomous agent for eda,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2024.
- [15] M. Liu, T.-D. Ene, R. Kirby, C. Cheng, N. Pinckney, R. Liang, J. Alben, H. Anand, S. Banerjee, I. Bayraktaroglu *et al.*, “Chipnemo: Domain-adapted llms for chip design,” *arXiv preprint arXiv:2311.00176*, 2023.
- [16] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [17] J. Blocklove *et al.*, “Chip-chat: Challenges and opportunities in conversational hardware design,” in *2023 ACM/IEEE 5th Workshop on Machine Learning for CAD (MLCAD)*. IEEE, Sep. 2023.
- [18] Y. Fu, Y. Zhang, Z. Yu, S. Li, Z. Ye, C. Li, C. Wan, and Y. C. Lin, “Gpt4aigchip: Towards next-generation ai accelerator design automation via large language models,” in *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. IEEE, 2023, pp. 1–9.
- [19] K. Chang, Y. Wang, H. Ren, M. Wang, S. Liang, Y. Han, H. Li, and X. Li, “Chipgpt: How far are we from natural language hardware design,” *arXiv preprint arXiv:2305.14019*, 2023.
- [20] Y. Yang, F. Teng, P. Liu, M. Qi, C. Lv, J. Li, X. Zhang, and Z. He, “Haven: Hallucination-mitigated llm for verilog code generation aligned with hdl engineers,” *arXiv preprint arXiv:2501.04908*, 2025.
- [21] M. Liu, Y.-D. Tsai, W. Zhou, and H. Ren, “Craftrtl: High-quality synthetic data generation for verilog code models with correct-by-construction non-textual representations and targeted code repair,” *arXiv preprint arXiv:2409.12993*, 2024.
- [22] Y. Liu, C. Xu, Y. Zhou, Z. Li, and Q. Xu, “Deeprtl: Bridging verilog understanding and generation with a unified representation model,” *arXiv preprint arXiv:2502.15832*, 2025.
- [23] Y. Zhao, H. Zhang, H. Huang, Z. Yu, and J. Zhao, “Mage: A multi-agent engine for automated rtl code generation,” *arXiv preprint arXiv:2412.07822*, 2024.
- [24] F. Cui, C. Yin, K. Zhou, Y. Xiao, G. Sun, Q. Xu, Q. Guo, D. Song, D. Lin, X. Zhang *et al.*, “Origen: Enhancing rtl code generation with code-to-code augmentation and self-reflection,” *arXiv preprint arXiv:2407.16237*, 2024.
- [25] H. Pearce, B. Ahmad, B. Tan, B. Dolan-Gavitt, and R. Karri, “Asleep at the keyboard? assessing the security of github copilot’s code contributions,” *Communications of the ACM*, vol. 68, no. 2, pp. 96–105, 2025.
- [26] Z. Ji, P. Ma, and S. Wang, “Unlearnable examples: Protecting open-source software from unauthorized neural code learning,” in *SEKE*, 2022, pp. 525–530.
- [27] Z. Yu, Y. Wu, N. Zhang, C. Wang, Y. Vorobeychik, and C. Xiao, “Codeiprompt: intellectual property infringement assessment of code language models,” in *International conference on machine learning*. PMLR, 2023, pp. 40 373–40 389.
- [28] H. Du, S. Liu, L. Zheng, Y. Cao, A. Nakamura, and L. Chen, “Privacy in fine-tuning large language models: Attacks, defenses, and future directions,” *arXiv preprint arXiv:2412.16504*, 2024.
- [29] N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson *et al.*, “Extracting training data from large language models,” in *30th USENIX security symposium (USENIX Security 21)*, 2021, pp. 2633–2650.
- [30] Z. Sun, X. Du, F. Song, M. Ni, and L. Li, “Coprotector: Protect open-source code against unauthorized training usage with data poisoning,” in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 652–660.
- [31] L. Niu, S. Mirza, Z. Maradni, and C. Pöpper, “{CodexLeaks}: Privacy leaks from code generation language models in {GitHub} copilot,” in *32nd USENIX Security Symposium (USENIX Security 23)*, 2023, pp. 2133–2150.
- [32] R. Schuster, C. Song, E. Tromer, and V. Shmatikov, “You autocomplete me: Poisoning vulnerabilities in neural code completion,” in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1559–1575.
- [33] H. Yang, K. Xiang, M. Ge, H. Li, R. Lu, and S. Yu, “A comprehensive overview of backdoor attacks in large language models within communication networks,” *IEEE Network*, 2024.
- [34] R. Liu, T. Wang, Y. Cao, and L. Xiong, “Precurious: How innocent pre-trained language models turn into privacy traps,” in *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security*, 2024, pp. 3511–3524.
- [35] M. S. Ozdayi, C. Peris, J. FitzGerald, C. Dupuy, J. Majmudar, H. Khan, R. Parikh, and R. Gupta, “Controlling the extraction of memorized data from large language models via prompt-tuning,” *arXiv preprint arXiv:2305.11759*, 2023.
- [36] F. He, T. Zhu, D. Ye, B. Liu, W. Zhou, and P. S. Yu, “The emerged security and privacy of llm agent: A survey with case studies,” *arXiv preprint arXiv:2407.19354*, 2024.
- [37] V. Rathod, S. Nabavirazavi, S. Zad, and S. S. Iyengar, “Privacy and security challenges in large language models,” in *2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC)*. IEEE, 2025, pp. 00 746–00 752.
- [38] L. L. Mankali, J. Bhandari, M. Alam, R. Karri, M. Maniatakos, O. Sinanoglu, and J. Knechtel, “Rtl-breaker: Assessing the security of llms against backdoor attacks on hdl code generation,” *arXiv preprint arXiv:2411.17569*, 2024.
- [39] S. Liu, Y. Yao, J. Jia, S. Casper, N. Baracaldo, P. Hase, Y. Yao, C. Y. Liu, X. Xu, H. Li *et al.*, “Rethinking machine unlearning for large language models,” *Nature Machine Intelligence*, pp. 1–14, 2025.

- [40] J. Ji, Y. Liu, Y. Zhang, G. Liu, R. Kompella, S. Liu, and S. Chang, “Reversing the forget-retain objectives: An efficient llm unlearning framework from logit difference,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 12 581–12 611, 2024.
- [41] C. Liu, Y. Wang, J. Flanigan, and Y. Liu, “Large language model unlearning via embedding-corrupted prompts,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 118 198–118 266, 2024.
- [42] M. Choi, D. Rim, D. Lee, and J. Choo, “Snap: Unlearning selective knowledge in large language models with negative instructions,” *arXiv preprint arXiv:2406.12329*, 2024.
- [43] Y. Wang, J. Wei, C. Y. Liu, J. Pang, Q. Liu, A. P. Shah, Y. Bao, Y. Liu, and W. Wei, “Llm unlearning via loss adjustment with only forget data,” *arXiv preprint arXiv:2410.11143*, 2024.
- [44] P. Maini, Z. Feng, A. Schwarzschild, Z. C. Lipton, and J. Z. Kolter, “Tofu: A task of fictitious unlearning for llms,” *arXiv preprint arXiv:2401.06121*, 2024.
- [45] R. Zhang, L. Lin, Y. Bai, and S. Mei, “Negative preference optimization: From catastrophic collapse to effective unlearning,” *arXiv preprint arXiv:2404.05868*, 2024.
- [46] C. Fan, J. Liu, L. Lin, J. Jia, R. Zhang, S. Mei, and S. Liu, “Simplicity prevails: Rethinking negative preference optimization for llm unlearning,” *arXiv preprint arXiv:2410.07163*, 2024.
- [47] N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, J. D. Li, A.-K. Dombrowski, S. Goel, L. Phan *et al.*, “The wmdp benchmark: Measuring and reducing malicious use with unlearning,” *arXiv preprint arXiv:2403.03218*, 2024.
- [48] A. Allam and M. Shalan, “Rtl-repo: A benchmark for evaluating llms on large-scale rtl design projects,” *arXiv preprint arXiv:2405.17378*, 2024.
- [49] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text summarization branches out*, 2004, pp. 74–81.
- [50] W. Shi, A. Ajith, M. Xia, Y. Huang, D. Liu, T. Blevins, D. Chen, and L. Zettlemoyer, “Detecting pretraining data from large language models,” *arXiv preprint arXiv:2310.16789*, 2023.
- [51] J. Zhang, J. Sun, E. Yeats, Y. Ouyang, M. Kuo, J. Zhang, H. F. Yang, and H. Li, “Min-k%+: Improved baseline for detecting pre-training data from large language models,” *arXiv preprint arXiv:2404.02936*, 2024.
- [52] V. B. Kumar, R. Gangadharaiah, and D. Roth, “Privacy adhering machine un-learning in nlp,” *arXiv preprint arXiv:2212.09573*, 2022.

APPENDIX

A. Mathematical Formulation of Unlearning

To formally model machine unlearning in the context of LLM-aided RTL generation, we define two disjoint datasets: the *Retain Dataset* \mathcal{D}_r and the *Forget Dataset* \mathcal{D}_f , where $\mathcal{D}_r \cap \mathcal{D}_f = \emptyset$. The goal of unlearning is to update a fine-tuned model f_θ such that it forgets knowledge gained from \mathcal{D}_f while preserving performance on \mathcal{D}_r . Mathematically, this is achieved by optimizing a modified objective function that induces high loss on \mathcal{D}_f and low loss on \mathcal{D}_r . A common formulation involves solving:

$$\min_{\theta} \mathcal{L}_{\text{retain}}(\theta) - \lambda \cdot \mathcal{L}_{\text{forget}}(\theta),$$

where $\mathcal{L}_{\text{retain}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_r}[\ell(f_\theta(x), y)]$ and $\mathcal{L}_{\text{forget}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_f}[\ell(f_\theta(x), y)]$, and $\lambda > 0$ controls the unlearning aggressiveness. Techniques like Gradient Ascent (GA) perform unlearning by maximizing $\mathcal{L}_{\text{forget}}$ through reversed gradients, while Gradient Difference (GD) introduces a difference term between gradients of \mathcal{D}_f and \mathcal{D}_r to encourage selective forgetting. Preference-based methods, such as Preference Optimization (PO) and Negative Preference Optimization (NPO), further refine this by optimizing the model’s alignment with desired or misdirected responses on \mathcal{D}_f , effectively pushing its latent representations away from those associated with proprietary or malicious logic. Such formulations enforce representational decoupling from the Forget Dataset, thereby

Algorithm 1 Machine Unlearning

Require: Retain dataset \mathcal{D}_r , Forget dataset \mathcal{D}_f
Initialize LLM f_θ with pre-trained weights
Define loss: $\mathcal{L}_{\text{unlearn}}(\theta) = \mathcal{L}_{\text{retain}}(\theta) - \lambda \cdot \mathcal{L}_{\text{forget}}(\theta)$
for $i = 1$ to EPOCHS **do**
 Compute $\nabla_{\theta} \mathcal{L}_{\text{unlearn}}$
 Update parameters: $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{unlearn}}$
end for

mitigating information leakage and contamination in downstream RTL synthesis tasks.

B. Unlearning Evaluation Metrics

- **Forget Probability (FP):** This metric evaluates each instance in the retain or forget set by computing the normalized token-level probability of the answer, reflecting how confidently the model predicts the answer given the question. Specifically, we calculate the following:

$$P(y | x)^{1/|y|}$$

where x denotes the question, y is the corresponding answer, and $|y|$ represents the number of tokens in the answer.

- **Forget ROUGE (FR):** This metric computes the ROUGE-L recall score between the ground truth answer y and the model-generated answer \hat{y} for each sample in the forget dataset \mathcal{D}_f . The ROUGE-L recall is defined as:

$$\text{ROUGE-L Recall} = \frac{LCS(y, \hat{y})}{|y|}$$

where $LCS(y, \hat{y})$ denotes the length of the longest common subsequence between y and \hat{y} , and $|y|$ is the length of the ground truth answer. Lower ROUGE-L recall scores indicate better unlearning performance.

- **Min-K% and Min-K%++:** The Min-K% metric focuses on the model’s confidence in generating tokens. For each generated sequence, it identifies the $k\%$ tokens with the lowest predicted probabilities and computes their average log-likelihood:

$$\text{Min-K\%} = \frac{1}{k} \sum_{i=1}^k \log p(t_i | x)$$

where t_i are the tokens in the bottom $k\%$ of predicted probabilities. Min-K%++ enhances this by calibrating based on token distribution statistics, providing a more robust detection of memorized content.

- **PrivLeak:** This metric assesses the privacy risk by measuring the difference in the Area Under the Receiver Operating Characteristic Curve (AUC-ROC) between the unlearned model f_{unlearn} and a retrained model f_{retrain} on the forget dataset \mathcal{D}_f . It is defined as:

$$\text{PrivLeak} = \text{AUC}_{f_{\text{unlearn}}} - \text{AUC}_{f_{\text{retrain}}}$$

A significant deviation from zero indicates a higher privacy risk, suggesting that the unlearned model retains information from the forget dataset.

C. Extra Results for Other Metrics

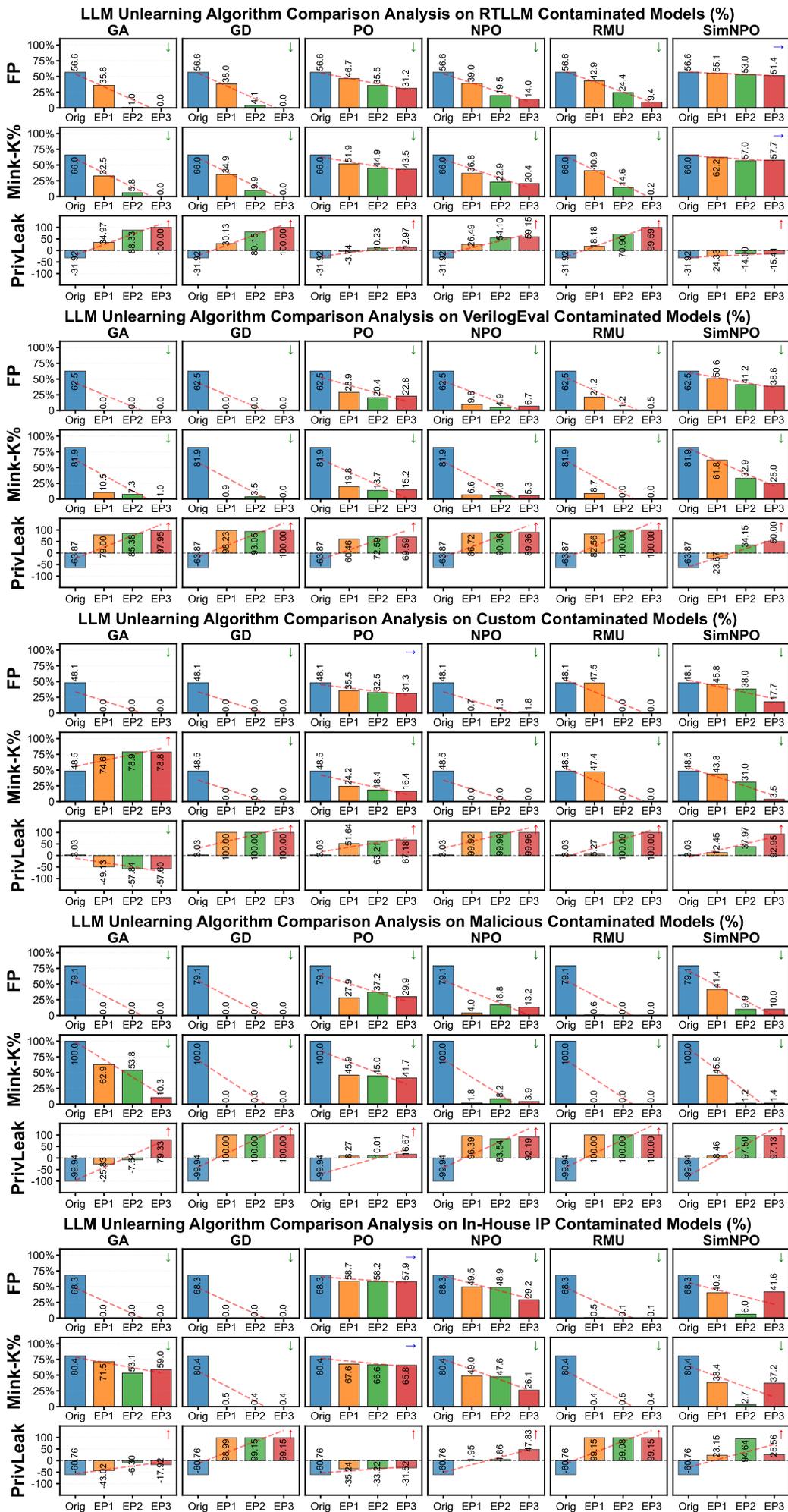


Fig. 7. Unlearn Contaminated Models with extra evaluation metrics