

# When Forgetting Builds Reliability: LLM Unlearning for Reliable Hardware Code Generation

Yiwen Liang<sup>†</sup>, Qiufeng Li<sup>†</sup>, Shikai Wang, Weidong Cao\*

Department of ECE, The George Washington University, Washington, D.C., USA

{yiwen.liang, qiufeng.li, shikai.wang, weidong.cao}@gwu.edu

<sup>†</sup>equal contribution, \*corresponding author

**Abstract**—Large Language Models (LLMs) have shown strong potential in accelerating digital hardware design through automated code generation. Yet, ensuring their reliability remains a critical challenge, as existing LLMs trained on massive heterogeneous datasets often exhibit problematic memorization of proprietary intellectual property (IP), contaminated benchmarks, and unsafe coding patterns. To mitigate these risks, we propose a novel unlearning framework tailored for LLM-based hardware code generation. Our method combines (i) a syntax-preserving unlearning strategy that safeguards the structural integrity of hardware code during forgetting, and (ii) a fine-grained floor-aware selective loss that enables precise and efficient removal of problematic knowledge. This integration achieves effective unlearning without degrading LLM code generation capabilities. Extensive experiments show that our framework supports forget sets up to  $3\times$  larger, typically requiring only a single training epoch, while preserving both syntactic correctness and functional integrity of register-transfer level (RTL) codes. Our work paves an avenue towards reliable LLM-assisted hardware design.

## I. INTRODUCTION

Recent advances in Large Language Models (LLMs) have ushered in a new paradigm for hardware development. General-purpose models such as GPT, and domain-specific models like RTLCode [1] and VeriGen [2], have demonstrated remarkable capabilities in generating hardware codes, thereby substantially accelerating digital design workflows, shortening development cycles, and reducing manual engineering efforts. Despite this progress, fundamental challenges remain to ensure the reliability of LLM-based code generation. Trained on massive and heterogeneous datasets, current LLMs often memorize undesirable content, including proprietary intellectual property (IP), contaminated benchmarks, and unsafe coding patterns, which compromise evaluation fairness, risk IP leakage, and threaten the reliability of downstream hardware systems [3]. For example, benchmark frameworks such as VerilogEval [4] and RtlLm [5] have been shown to suffer from extensive contamination. Recent studies even reveal that up to 40% of LLM-generated code contains security vulnerabilities [6]. Moreover, ongoing lawsuits alleging unauthorized reproduction of proprietary code by LLMs further underscore these risks [7]. Together, these challenges raise a pressing question: How can we intentionally mitigate the influence of problematic LLM memorization to achieve reliable hardware code generation?

Existing efforts to mitigate problematic memorization in general LLMs have primarily focused on retraining and machine unlearning. Retraining from scratch on sanitized datasets represents the most reliable way to achieve exact unlearning; yet, the prohibitive resource and time costs of training large-scale models make this approach largely impractical in real-world settings. For instance, training the LLaMA 3.1 model

(8B parameters) [8] from scratch requires approximately 1.46 million GPU hours on NVIDIA H100-80GB GPUs. Machine unlearning has emerged as a more practical alternative [9]–[15], which aims to selectively remove the influence of specific training data while preserving the utility of models. This ensures that the resulting model behaves as if the removed data had never been included in training. Specifically, existing efforts adopt post-training strategies for efficient unlearning without retraining [9]–[15], and have shown promise in natural language processing tasks, including harmful response removal, copyright erasure, and privacy protection. Yet, the application of unlearning to hardware code generation remains significantly underdeveloped. Recently, only SALAD [16] applied existing unlearning techniques to hardware design without any domain adaptation, and therefore suffers from substantial utility degradation. This gap highlights the need to develop domain-specific unlearning strategies to ensure reliable hardware code generation without compromising utility.

This paper proposes the first-of-its-kind domain-specific unlearning framework tailored to reliable hardware code generation, which is capable of eliminating problematic knowledge while maintaining reliable hardware code synthesis. Specifically, by recognizing that RTL codes require strict syntax, precise dependency modeling, and faithful execution semantics, our strategies incorporate these factors into unlearning processes to strike a balance between forgetting effectiveness and hardware code generation reliability. Our key contributions are as follows:

- **Novel unlearning framework for reliable hardware code generation.** We pioneer the development of domain-specific unlearning frameworks to significantly improve the reliability of LLM-based hardware code generation with tailored strategies.
- **Syntax-preserving unlearning strategy.** We introduce an RTL-specific mechanism that selectively excludes structural tokens from the forgetting process, ensuring that unlearning does not compromise the syntactic integrity of generated hardware code.
- **Fine-grained floor-aware selective loss (FiFSL).** We develop a novel token-level objective that combines margin-based control with selective averaging, enabling precise and accelerated forgetting of undesired samples.
- **Scalable and robust unlearning for hardware code generation.** Comprehensive experiments show that our framework achieves up to  $3\times$  forget sets than traditional methods and requires only a single training epoch, while maintaining syntactic validity and functional reliability

in RTL code generation, demonstrating state-of-the-art improvements in both efficiency, reliability, and utility.

## II. BACKGROUND

### A. Syntax and Semantics of Hardware Description Language

A typical digital module is synchronous, consisting of two interacting components: **combinational logic**, which computes functions of the current inputs without memory, and **sequential logic**, which stores and updates states through registers or flip-flops at each clock edge. Register-Transfer Level (RTL) modeling captures this behavior by specifying data transfers between registers under clock control, together with the combinational logic that generates next-state values. In Verilog, for instance, sequential logic is commonly expressed using `always @(posedge clk or negedge rst_n)` blocks with non-blocking assignments (`<=`), while combinational logic is described with `assign` statements or `always @(*)` blocks using blocking assignments (`=`). Consequently, Verilog code for digital designs frequently relies on recurring syntax patterns, particularly `always` blocks that encode both combinational and sequential behavior.

### B. Hardware Design with LLM

Generative AI is rapidly reshaping the landscape of hardware development [17]–[22]. In particular, the integration of Large Language Models (LLMs) into digital hardware design has shown remarkable progress, demonstrating strong capability in Verilog generation, assertion synthesis, testbench creation, and optimization of electronic design automation (EDA) workflows [23]–[26]. Early efforts such as ChipNemo [26] and ChipGPT [27], show that fine-tuning and prompt engineering can substantially improve RTL generation quality. More recent frameworks like HAVEN [28], CraftRTL [29], and MAGE [30] extend these capabilities through non-textual representations and multi-agent strategies. Complementary benchmarks such as VerilogEval [4] and RtlM [5] further confirm that LLMs can produce syntactically correct and functionally reliable RTL, highlighting their promise to accelerate design workflows.

Despite these advances, fundamental challenges remain in ensuring the safety, legality, and reliability of LLM-based hardware code generation. Existing LLMs trained on massive public repositories risk propagating vulnerabilities and insecure coding practices through the memorization of buggy or outdated patterns. To address these risks, security-focused techniques have emerged, such as SafeCoder [6], which augments training with curated secure-code datasets, and VersiCode [31], a benchmark for evaluating robustness to Application Programming Interface (API) changes across releases. Yet, these approaches fall short of providing comprehensive guarantees. SafeCoder depends heavily on dataset quality and coverage, which cannot anticipate all insecure coding patterns or future vulnerabilities, leaving models susceptible to unseen exploits. VersiCode, by contrast, is limited to version-specific API robustness and does not address broader issues such as memorization of proprietary code, syntactic validity in hardware description languages, or broader security vulnerabilities. These limitations highlight the

need for novel approaches that can suppress unsafe behaviors while preserving the utility of LLMs for code generation.

### C. Machine Unlearning

Machine unlearning, the process of eliminating the influence of specific data from a trained model without retraining from scratch, has recently emerged as a critical research field. In the context of LLMs, unlearning has been increasingly explored to address applications such as privacy preservation, removal of copyrighted material, and suppression of harmful content [14]. Early methods, such as Gradient Ascent (GA) [9], attempt to “reverse” prior learning by performing ascent on the next-token prediction loss of the forget set, but they often suffer from instability and collateral drift on the retain distribution due to the indiscriminate and unbounded optimization signal. To tackle these issues, Negative Preference Optimization (NPO) [11] introduces an alignment-inspired loss function embedded with a smooth bounded penalty to promote divergence from forget-set behavior while maintaining stability. More recently, Sim-NPO [12] refines this approach by removing the reliance on a teacher model and adopting a reference-free, length-normalized objective, thereby mitigating reference bias and enabling more balanced forgetting with stronger utility preservation.

While these advances have demonstrated the feasibility of unlearning in natural language processing (NLP), their applicability to hardware code generation remains largely unexplored. Unlike natural language, RTL code exhibits strict syntactic and structural constraints, where even minor inconsistencies can render outputs unsynthesizable. As a result, NLP-oriented unlearning methods cannot be directly transferred to hardware description languages—as shown by a recent work SALAD [16], ignoring domain-specific syntax significantly undermines the validity of code generation. These limitations highlight the need for domain-specific unlearning strategies that not only suppress undesired behaviors but also preserve the correctness and reliability of RTL code generation for hardware design.

## III. METHODOLOGY

### A. Framework Overview

We propose a domain-specific unlearning framework for reliable LLM-based hardware code generation. Fig. 1 shows its overview, which leverages harmful data as a forget set, enabling the removal of problematic knowledge while preserving code utility. Our approach combines two major techniques: (1) **syntax-preserving masking**, which protects reserved keywords and code spans to maintain compilability, and (2) **fine-grained floor-aware selective loss (FiFSL)**, which applies margin-based forgetting pressure while gating out already forgotten samples to focus updates on the hardest cases. Together, these strategies enable efficient, stable, and utility-preserving unlearning, yielding LLMs that generate reliable hardware codes.

### B. Problem Definition

We consider the problem of machine unlearning in the context of LLMs for RTL code generation. Let  $\pi_\theta$  denote a pretrained LLM parameterized by  $\theta$ , which generates RTL code  $y$  given an instruction or design intent  $x$ , i.e.,  $\pi_\theta(y|x)$ .

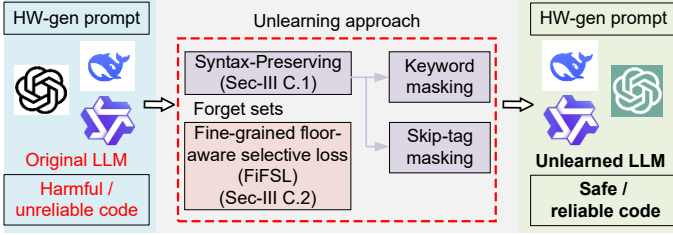


Fig. 1. Overview of our proposed domain-specific LLM unlearning framework for reliable hardware code generation (HW-gen).

In practice, not all training data should be retained indefinitely. Certain samples may correspond to privacy-sensitive hardware designs, erroneous code snippets, or legally restricted RTL modules [4]. Such cases can arise due to incomplete data cleaning, outdated versioned code, etc. We denote the set of data to be forgotten as  $D_f = (x_f, y_f)$ . The objective of unlearning is to update the model into  $\pi_\theta^*$  such that its dependence on  $D_f$  is minimized. Specifically, given a prompt  $x_f$ , the outputs of  $\pi_\theta^*$  should diverge from the undesired targets  $y_f$ . Meanwhile, the model must preserve its utility on the remaining distribution. We denote the validation set by  $D_v = (x_v, y_v)$ , which reflects the task-relevant code generation capability to be preserved. Formally, the unlearning goal can be expressed as:

$$\text{Find } \pi_\theta^* \text{ s.t. } \pi_\theta^*(y_f | x_f) \text{ is minimized, } \forall (x_f, y_f) \in D_f, \\ \pi_\theta^*(y_v | x_v) \approx \pi_\theta(y_v | x_v), \forall (x_v, y_v) \in D_v. \quad (1)$$

This formulation captures the dual goal of unlearning: removing the influence of undesired data while retaining the model's competence in hardware code generation. To address this dual goal in practice, our methodology consists of two key components: (i) syntax-preserving unlearning that ensures syntax structural correctness is not disrupted, and (ii) a fine-grained floor-aware selective loss function that accelerates and strengthens the forgetting process.

### C. Proposed Unlearning Framework

1) *Syntax-Preserving Unlearning*: Unlike free-form natural languages, RTL codes follow a finite and well-defined grammar. Their tokens, such as reserved keywords, operators, delimiters, and numeric literals, recur systematically to capture combinatorial and sequential behaviors. While exact idioms differ across cases (e.g., continuous assignments or sensitivity lists for combinational logic vs. clocked processes for sequential logic), such tokens represent universal syntactic backbones rather than task-specific semantics. Thus, existing unlearning strategies [16] that indiscriminately suppress these high-frequency syntax tokens are harmful: they undermine grammaticality and compromise the compilability and synthesizability of RTL programs (i.e., treating every token in the toy example Verilog code of a 4-bit counter (Fig. 2(a)) as removable). Effective unlearning must be syntax-preserving, targeting semantic content while enforcing grammar constraints to maintain well-formed hardware description languages (HDLs).

To achieve this, we propose a **syntax-preserving unlearning strategy that selectively protects structural elements while forgetting only task-specific tokens**. As shown in Fig. 2(b),

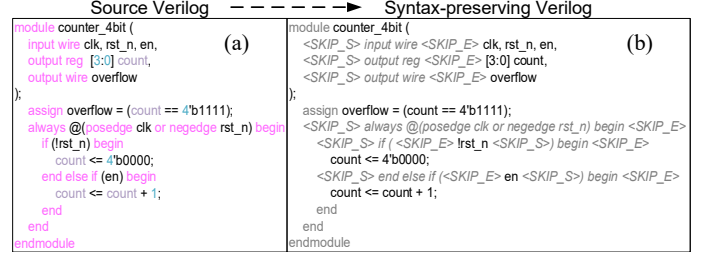


Fig. 2. Syntax-preserving unlearning: source Verilog (a) vs. syntax-preserving Verilog (b). The preserved syntax is masked in gray.

our approach integrates two complementary mechanisms to preserve structural correctness, grounded in the IEEE hardware code standard [32]–[34], and excludes them from the loss function. First, keyword masking is applied to a curated set of Verilog reserved words such as *module*, *assign*, *wire*, etc, ensuring that these essential tokens are not modified by forgetting operations. Second, skip-tag masking introduces special tokens  $\langle \text{SKIP\_S} \rangle$  and  $\langle \text{SKIP\_E} \rangle$  to mark spans of texts that should be ignored during loss computation, and all tokens appearing between these tags are excluded from the unlearning objective. To unify both operations, we assign each token  $x_t$  in an input sequence  $\mathbf{x} = (x_1, \dots, x_T)$  with a binary mask value  $m_t$ , where  $m_t = 0$  if  $x_t$  falls within a skip-tag span or belongs to the reserved keyword set  $\mathcal{K}$ , and  $m_t = 1$  otherwise. The loss for syntax-preserving unlearning is thus computed directly as

$$\mathcal{L}_{\text{syntax}} = \sum_{t=1}^T m_t \ell(\pi_\theta(x_t | \mathbf{x}_{<t}), y_t) \quad (2)$$

with  $\ell(\cdot)$  denoting the token-level cross-entropy. By aligning the masking rule with the loss function, the model forgets undesired content while maintaining the universal syntactic backbone required for compilable RTL code. It is important to note that from the computational perspective, the proposed masking strategy introduces negligible overheads, since it only requires constructing a token-level binary mask during pre-processing and applying it element-wise to the loss function. This makes syntax-preserving unlearning highly efficient and scalable, enabling its deployment in large-scale Verilog LLMs without incurring noticeable computational or memory costs.

2) *FiFSL: Fine-grained Floor-aware Selective Loss*: To further enhance the unlearning process, we introduce a **fine-grained floor-aware selective loss (FiFSL)** that operates at the token level and combines margin-based control with a selective averaging mechanism. The key idea is to impose a lower bound on the loss to prevent excessive forgetting, while gradients are computed only for unmasked tokens and unforbidden samples.

**Forward**: In the forward pass, for each sample  $i$  in a batch of size  $b$ , let  $\mathcal{V}_i$  denote the set of all valid (unmasked) tokens after the syntax-preserving masking in Sec. III-C1. The token-normalized negative log-likelihood for sample  $i$  is defined as

$$L_i = 1/|\mathcal{V}_i| \cdot \mathcal{L}_{\text{syntax}} \quad (3)$$

To impose forgetting, we shift this loss by a margin parameter  $\gamma$  and apply a smooth negative-preference mapping, yielding

$$\phi_i = 2/\beta \text{softplus}(-\beta(L_i - \gamma)), \quad (4)$$

where  $\beta > 0$  controls the curvature and saturation rate (larger  $\beta$  gives a sharper transition), and  $\text{softplus}(z) = \log(1 + e^z)$  is used to provide a smooth, bounded mapping. This construction exerts strong forgetting pressure when a sample is under-forgotten ( $L_i < \gamma$ ), and gradually saturates once the margin is exceeded, thus preventing runaway gradients. To further prevent excessive forgetting and suppress outlier effects, we introduce a hard floor  $L_{\min}$  and average only over active samples whose penalties exceed this threshold ( $\phi_i > L_{\min}$ ), contributing both to the loss and to the gradient. Let  $\alpha_i = \mathbf{1}\{\phi_i > L_{\min}\}$  and  $N_{\text{act}} = \sum_i \alpha_i$ , the mini-batch objective loss is expressed as

$$\mathcal{L}_{\text{FiFSL}}^{\text{batch}} = 1/N_{\text{act}} \sum_i \alpha_i \phi_i. \quad (5)$$

This is algebraically equivalent to first clamping  $\phi_i$  by  $L_{\min}$ , masking inactive samples, and then normalizing by the number of active ones. As a result, samples with  $\phi_i \leq L_{\min}$  contribute no gradient and are effectively excluded from further updates. At the dataset level, FiFSL can be expressed as the expectation over the forget set  $D_f$ :

$$\begin{aligned} \mathcal{L}_{\text{FiFSL}}(\theta) &= \mathbb{E}(x, y) \in D_f \left[ \alpha(x, y) \cdot \phi(\pi_\theta, x, y) \right], \\ \alpha(x, y) &= \mathbf{1}\{\phi(\pi_\theta, x, y) > L_{\min}\}. \end{aligned} \quad (6)$$

**Backpropagation:** By establishing the forward objective, we now analyze the backward signal to show how our loss enforces forgetting stably and selectively. Recall Eqs. (3) and (4). Using  $\text{softplus}'(z) = \sigma(z)$  with the  $\sigma$  as logistic sigmoid, differentiating  $\phi_i$  with respect to the per-sample loss yields

$$\partial \phi_i / \partial L_i = -2\sigma(-\beta(L_i - \gamma)). \quad (7)$$

Thus, the sensitivity lies strictly between  $-2$  and  $0$ : when  $L_i < \gamma$ , the term  $-\beta(L_i - \gamma)$  is positive,  $\sigma(\cdot) \approx 1$ , and the gradient magnitude approaches two, producing a strong negative push that increases  $L_i$  and drives forgetting. Consequently, when  $L_i > \gamma$ , the sigmoid gradually reduces the gradient, leading to milder updates and preventing overly large gradients that could cause excessive forgetting. At the token level, the derivative of  $L$  with respect to a logit at position  $t$  is  $\frac{1}{|\mathcal{V}_i|} m_{i,t} (\mathbf{p}_{i,t} - \mathbf{e}(y_{i,t}))$ . Here,  $m_{i,t}$  is the syntax mask,  $\mathbf{p}_{i,t}$  the softmax vector, and  $\mathbf{e}(y_{i,t})$  the one-hot target. The normalization by  $|\mathcal{V}_i|$  ensures that gradients are comparable across sequences of different lengths, while the mask guarantees that reserved keywords or skipped tokens never receive updates. Combining these results, the batch-level gradient becomes

$$\frac{\partial \mathcal{L}_{\text{FiFSL}}^{\text{batch}}}{\partial \text{logit}_{i,t}} = \frac{\alpha_i}{N_{\text{act}}} (-2\sigma(-\beta(L_i - \gamma))) \frac{1}{|\mathcal{V}_i|} m_{i,t} (\mathbf{p}_{i,t} - \mathbf{e}(y_{i,t})). \quad (8)$$

This expression shows that FiFSL rescales the standard cross-entropy gradient by a bounded negative factor and suppresses it entirely when either the token is masked or the sample is already forgotten ( $\alpha_i = 0$ ).

The implications are twofold. First, FiFSL enforces forgetting with control: active samples below the margin receive strong forgetting pressure, but this pressure decays smoothly once the margin is passed, stabilizing optimization. Second, FiFSL avoids unnecessary updates by gating entire samples

---

**Algorithm 1:** Our domain-specific unlearning method.

---

```

Input: Forget set  $\mathcal{D}_f$ , model  $\pi_\theta$ , reserved keywords  $K$ ,
skip-tags
foreach batch  $(x, y) \in \mathcal{D}_f$  do
  // Syntax-preserving
  foreach token  $x_t$  in  $x$  do
    if  $x_t \in K$  or inside skip-tags then
       $m_t \leftarrow 0$ ;
       $\text{labels}[t] \leftarrow -100$  // ignored in loss
    else
       $m_t \leftarrow 1$ 
   $L_{\text{syntax}} \leftarrow \sum_t m_t \cdot \ell(\pi_\theta(x_t|x_{<t}), y_t)$ ;
  // FiFSL
  foreach sample  $i$  do
     $L_i \leftarrow \text{normalized}(L_{\text{syntax}}(i))$ ;
     $\phi_i \leftarrow \frac{2}{\beta} \text{softplus}(-\beta(L_i - \gamma))$ ;
     $\alpha_i \leftarrow 1$  if  $\phi_i > L_{\min}$  else  $0$ ;
   $L_{\text{batch}} \leftarrow \frac{\sum_i \alpha_i \phi_i}{\sum_i \alpha_i}$ ;
   $\theta \leftarrow \theta - \eta \nabla_\theta L_{\text{batch}}$ ;

```

---

TABLE I  
PERFORMANCE OF BENCHMARK MODELS AFTER FINE-TUNING.

Model	Forget quality		Model utility			
	PrivLeak	MinK++	Loss	Pass@1	BLEU	chrF
Llama-7B	0.98	0.80	0.30	49%	50.7	59.1
Llama-8B	0.96	0.97	0.41	53%	40.1	56.2
DeepSeek-7B	0.97	0.83	0.28	55%	48.5	57.9
Qwen-7B	0.87	0.96	0.25	57%	50	57.7

once they have been sufficiently forgotten, unlike conventional unlearning techniques, which continue to update all samples uniformly. In practice, as training proceeds, the number of active samples decreases, and optimization effort concentrates on the hard-to-forget tails of the distribution. This selective mechanism not only prevents collateral degradation of syntax tokens but also reduces the number of epochs required to achieve effective unlearning compared to conventional methods.

3) *Summary of the Proposed Strategy:* In summary, our framework combines *syntax-preserving unlearning*, which safeguards the structural correctness of RTL code, with *FiFSL*, a margin-based selective forgetting loss that ensures stable and efficient removal of undesired knowledge. This joint design enables effective forgetting while preserving syntactic and functional validity. The process is presented in Algorithm 1.

## IV. EXPERIMENTS

### A. Experimental Methodology

1) *Evaluation Metrics:* Evaluations are conducted along two dimensions: (i) **forget quality** and (ii) **model utility**. These metrics are widely adopted in the LLM unlearning domain and are adapted here to the context of hardware code generation.

**Forget quality.** We employ two metrics to quantify the degree of forgetting: (1) *PrivLeak* [35], which quantifies the probability that the model regenerates HDL snippets or modules from the forget set. This corresponds to verbatim or near-verbatim leakage of RTL blocks (e.g., always blocks, state machines, or functional units). *Lower values indicate stronger forgetting and reduced IP leakage risk.* (2) *MinK++* [36], which evaluates the average log-likelihood of the lowest- $k\%$  gold tokens in generated HDL codes. By calibrating against

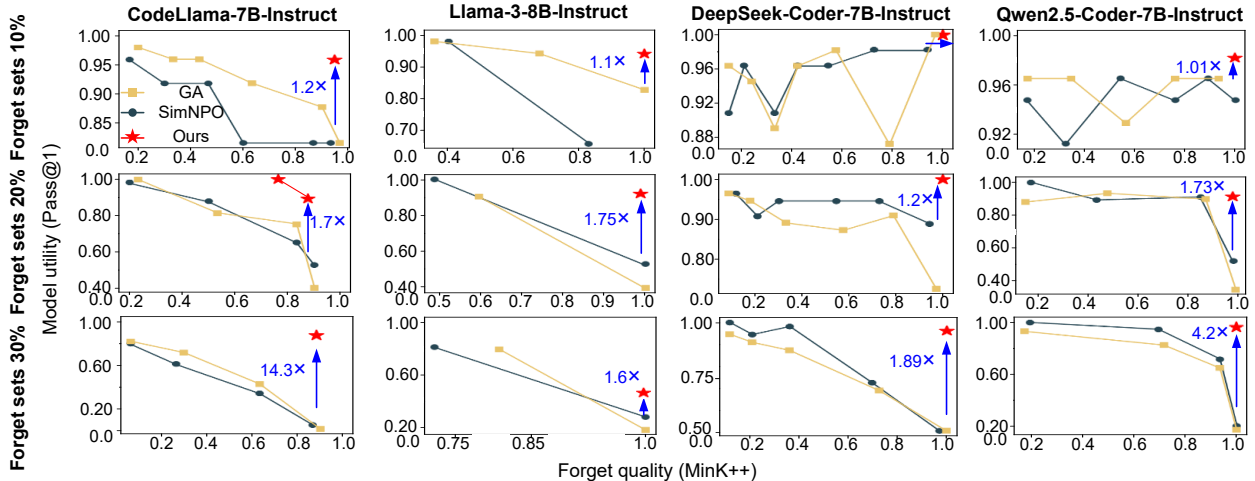


Fig. 3. Comparisons of forgetting quality and model utility across unlearning methods, varying forget-set sizes, and different LLMs. Ideal unlearning target: (1.0, 1.0), perfect forgetting with no utility loss. Our method achieves comparable or better trade-offs in only 1 epoch, while baselines require 4–7 epochs.

token distributions, it is sensitive to residual memorization of hardware coding patterns (e.g., signal naming conventions, pipeline templates) even when direct leakage does not occur. Lower scores suggest weaker memorization. For reference, untrained models typically yield scores around 0.5 [36]. **Model utility.** To ensure that forgetting does not degrade downstream performance, we evaluate model utility on hardware code generation tasks using: (1) *Cross-entropy loss* on held-out HDL test sets, measuring token-level predictive quality. (2) *Pass@1 accuracy* on functional hardware code generation benchmarks, assessing whether the generated Verilog/HDL compiles and passes simulation testbenches. (3) *BLEU* and *chrF*, which capture syntactic and stylistic fidelity. BLEU measures  $n$ -gram overlap with reference RTL implementations, while chrF evaluates character-level similarity, useful for ensuring correct use of hardware-specific syntax.

2) *Benchmarks and Experiment Setup:* We evaluate our approach on four LLMs fine-tuned for RTL code generation: Llama-3-8B-Instruct [8], CodeLlama-7B-Instruct [37], DeepSeek-Coder-7B-Instruct [38], and Qwen2.5-Coder-7B-Instruct [39]. Among them, Llama-3-8B-Instruct is a general-purpose model, while others are explicitly optimized for code understanding and generation. As the fine-tuned corpus, we adopt the open-source RTLCoder dataset [1], which provides instruction-response pairs comprising natural language design specifications and their corresponding RTL code implementations. Each model is fine-tuned on 1,000 samples for six epochs using the Adam optimizer with a learning rate of  $1 \times 10^{-5}$ , a maximum sequence length of 2,048 tokens, and a batch size of 2 on 4×A6000 GPUs. For inference, we set the sampling temperature to 0.5 and top- $p$  to 0.9. Table I reports the performance of these benchmark models after fine-tuning. All models exhibit strong memorization tendencies (PrivLeak and MinK++ > 0.8), which provides a baseline reference for subsequent unlearning evaluations.

Since there are no publicly available domain-specific datasets for proprietary or contaminated RTL designs, we follow prior works [9]–[16] to emulate realistic unlearning scenarios by subsampling 100, 200, and 300 examples from the training

corpus (i.e., about 10%, 20%, and 30%). These forget sets emulate practical cases where LLMs may memorize proprietary IP, contaminated benchmarks, or unsafe coding patterns that must be forgotten. The downstream utility is assessed on the held-out RTLCoder test set, consistent with the benchmark evaluation protocol.

3) *Unlearning Baselines and Configurations:* We compare our method against two of the most recent and representative LLM unlearning approaches, which SALAD [16] used. **Gradient Ascent (GA)** [9] maximizes the cross-entropy loss on forget samples:

$$\mathcal{L}_{\text{GA}}(\theta) = -\mathbb{E}_{(x,y) \in \mathcal{D}_f} [-\log \pi_{\theta}(y | x)]. \quad (9)$$

**SimNPO** [12] introduces a smooth penalty function to improve stability and avoid collapse of GA:

$$\mathcal{L}_{\text{SimNPO}}(\theta) = \mathbb{E}_{(x,y) \in \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\frac{\beta}{|y|} \log \pi_{\theta}(y | x) - \gamma \right) \right]. \quad (10)$$

To ensure fairness, we follow the reward-margin settings in [12] to set  $\beta = 2.5$  and  $\gamma = 0$ . Our **FIFSL** further introduces an adaptive floor parameter ( $L_{\min}$ ), chosen as 0.35, which corresponds to the validation loss after the first fine-tuning epoch. In addition, we explore a broader hyperparameter space by using Llama-7B as an example, varying  $\beta$  between 1.5 and 4 and  $\gamma$  from 0 to 2.5. Consistent with the findings of [12], the configuration  $\beta = 2.5, \gamma = 0$  achieves the best performance. We build our syntax-preserving dataset guided by [32]–[34]. Unlearning is performed using a learning rate of  $2 \times 10^{-6}$  until MinK++ converges near 0.5, ensuring effective forgetting without excessive utility degradation.

## B. Experimental Results

1) *Forgetting Effectiveness vs. Utility Preservation:* We compare our approach against GA and SimNPO in terms of forget quality and model utility. Fig. 3 shows results across different forget sets and benchmark models. We adopt MinK++ as the indicator of forget quality and Pass@1 as the measure of model utility, since they directly capture the trade-off between eliminating memorized knowledge and retaining downstream

TABLE II  
A COMPREHENSIVE COMPARISON OF MODEL PERFORMANCE AFTER UNLEARNING ON DIFFERENT SIZES OF FORGET SETS.

Model	Forget sets	Method	Epoch ↓	PrivLeak ↓	MinK++ ↓	Validation				Generalization Ability			
						Loss ↓	Pass@1 ↑	BLEU ↑	chrF ↑	Loss ↓	Pass@1 ↑	BLEU ↑	chrF ↑
Llama-7B	10%	GA	6	0.75	0.52	0.46	40%	40.5	50.3	0.70	68%	42.5	61.6
		SimNPO	6	0.75	0.51	0.46	40%	41.6	51.2	0.70	62%	40.6	61.1
		Ours	1	0.70	0.53	0.30	47%(-2%)	50.6	57.9	0.34	73%	46.0	67.0
	20%	GA	4	0.82	0.53	0.48	26%	40.5	48.8	0.67	53%	31.5	57.0
		SimNPO	4	0.82	0.53	0.50	20%	40.5	48.0	0.73	48%	30.0	55.8
		Ours	2	0.65	0.54	0.33	45%(-4%)	48.8	56.4	0.43	69%	35.3	59.1
	30%	GA	4	0.73	0.54	1.00	3%	13.2	28.7	1.70	5%	8.2	26.6
		SimNPO	4	0.72	0.53	1.00	1%	11.7	26.9	1.70	4%	8.0	24.7
		Ours	1	0.65	0.53	0.32	44%(-5%)	49.7	56.8	0.40	72%	37.9	60.0
Llama-8B	10%	GA	2	0.66	0.58	0.50	35%	35.0	50.2	0.70	45%	21.7	47.0
		SimNPO	3	0.61	0.50	0.51	44%	37.8	51.1	0.80	50%	21.7	47
		Ours	1	0.68	0.49	0.31	50%(-3%)	47.0	55.6	0.33	73%	44.2	62.4
	20%	GA	2	0.67	0.50	0.66	28%	30.3	44.2	1.10	30%	17.0	42.1
		SimNPO	2	0.69	0.50	0.59	21%	24.7	44.9	0.97	25%	12.5	37.0
		Ours	1	0.64	0.51	0.34	49%(-4%)	44.0	51.8	0.40	71%	41.9	58.3
	30%	GA	2	0.65	0.50	0.77	15%	16.8	36.9	1.20	24%	8.0	27.6
		SimNPO	2	0.65	0.50	0.78	10%	16.0	37.3	1.28	18%	7.5	27.5
		Ours	1	0.60	0.51	0.37	24%	35.0	50.0	0.53	34%	21.8	50.5
DeepSeek-7B	10%	GA	7	0.77	0.52	0.35	54%	44.2	53.3	0.42	81%	44.4	61.9
		SimNPO	7	0.75	0.50	0.37	55%	43.4	53.7	0.43	79%	45.1	61.7
		Ours	1	0.56	0.47	0.28	55%(+0%)	49.6	58.3	0.29	79%	36.7	59.3
	20%	GA	6	0.73	0.52	0.50	49%	31.5	48.6	0.73	79%	27.0	53.1
		SimNPO	6	0.72	0.51	0.57	40%	28.4	47.0	0.77	77%	24.3	51.5
		Ours	1	0.56	0.47	0.28	57%(+2%)	48.5	58.3	0.30	80%	37.6	58.2
	30%	GA	5	0.71	0.51	1.00	28%	19.2	42.3	1.10	16%	10.9	35.4
		SimNPO	5	0.71	0.50	0.73	28%	17.9	40.9	1.11	17%	10.8	34.8
		Ours	1	0.55	0.45	0.29	53%(-2%)	47.9	56.6	0.29	86%	41.1	62.0
Qwen-7B	10%	GA	6	0.51	0.47	0.37	54%	37.5	48.0	0.39	80%	46.3	58.1
		SimNPO	5	0.54	0.53	0.33	55%	41.0	50.0	0.39	81%	47.2	59.3
		Ours	1	0.60	0.47	0.25	56%(-1%)	49.2	57.3	0.28	85%	48.0	64.8
	20%	GA	4	0.61	0.51	0.43	30%	37.8	47.4	0.77	55%	29.4	50.9
		SimNPO	4	0.60	0.50	0.45	20%	36.4	45.9	0.79	42%	28.8	50.1
		Ours	1	0.57	0.51	0.25	53%(-4%)	46.9	54.7	0.32	83%	46.6	61.7
	30%	GA	3	0.64	0.50	0.49	13%	31.8	44.2	0.80	40%	24.0	41.0
		SimNPO	3	0.64	0.50	0.50	10%	34.9	44.8	0.88	45%	23.4	40.1
		Ours	1	0.56	0.49	0.27	55%(-2%)	46.2	55.0	0.34	86%	49.2	63.0

RTL code generation reliability. Both values are normalized for comparability. An ideal unlearning method aspires to the point (1.0, 1.0) in the forget-utility space, denoting perfect forgetting without utility loss. Our method consistently lies closest to this ultimate goal across all models and forget-set sizes. In particular, it enables forgetting up to 30% of the training corpus across coder LLMs, while conventional methods fail beyond 10% (i.e., a 3× larger forget set). For Llama-3-8B, our method achieves up to 20% forgetting, which we attribute to model-specific architectural constraints, as this model is not explicitly optimized for code understanding and generation. Nevertheless, this still outperforms conventional methods, demonstrating the robustness and broad applicability of our approach. Importantly, our method preserves the highest model utility among all methods: at 30% forgetting, it retains up to 14.3× higher Pass@1 than baselines, whereas GA and SimNPO exhibit sharp utility degradation. Concrete values across all metrics are reported in Table II. Moreover, our method accomplishes near-complete forgetting within a *single epoch*, drastically reducing unlearning cost, while GA and SimNPO require 4~7 epochs and underperform, particularly on larger forget sets.

2) *Generalization Ability and Scalability*: Table II reports all evaluation metrics, where “epoch” denotes the number of training epochs required to reach ideal forget quality, with ↓ (↑) indicating that lower (higher) values are better. Our method surpasses in all metrics. We further assess generalization on unseen, simple code generation tasks, shown as generalization ability columns. The shaded rows highlight that

TABLE III  
OVERALL COMPARISON WITH BASELINES.

Method	Forget Sets	Epochs	Utility
GA	1×	4~7	×
SimNPO	1×	4~7	×
Ours	2×~3×	1	✓

the unlearned models maintain Pass@1 close to their fine-tuned baselines, confirming that forgetting does not compromise general-purpose generation ability. In addition, the proposed approach sustains high BLEU and chrF scores across these tasks, demonstrating that effective unlearning can be achieved without sacrificing broad code generation performance.

3) *Summarization*: Overall, our domain-specific method supports up to 3× larger forget sets than existing methods (i.e., GA and SimNPO used by SALAD [16] without domain adaptation), achieves effective forgetting in only one epoch, and sustains both validation and generalization performance (Table III). These results show that it offers near-ideal unlearning: strong forgetting quality, preserved model utility, broad generalization, and exceptional efficiency, making it a scalable and practical solution for reliable code-specialized LLMs.

## V. CONCLUSION

We present the first domain-specific unlearning framework for RTL code generation, combining a syntax-preserving strategy with the FiFSL objective to enable precise and efficient forgetting while maintaining syntactic validity and functional reliability. Our method supports up to 3× larger forget sets with only one training epoch. Additional experiments are included in the updated arXiv version [40] to address reviewer concerns.

## REFERENCES

- [1] S. Liu, W. Fang, Y. Lu, J. Wang, Q. Zhang, H. Zhang, and Z. Xie, "Rtlcoder: Fully open-source and efficient llm-assisted rtl code generation technique," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2024.
- [2] S. Thakur, B. Ahmad, H. Pearce, B. Tan, B. Dolan-Gavitt, R. Karri, and S. Garg, "Verigen: A large language model for verilog code generation," *ACM Transactions on Design Automation of Electronic Systems*, vol. 29, no. 3, pp. 1–31, 2024.
- [3] Y. Dong, X. Jiang, H. Liu, Z. Jin, B. Gu, M. Yang, and G. Li, "Generalization or memorization: Data contamination and trustworthy evaluation for large language models," *arXiv preprint arXiv:2402.15938*, 2024.
- [4] 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.
- [5] Y. Lu, S. Liu, Q. Zhang, and Z. Xie, "Rtllm: 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.
- [6] J. He, M. Vero, G. Krasnopolska, and M. Vechev, "Instruction tuning for secure code generation," *arXiv preprint arXiv:2402.09497*, 2024.
- [7] D. Zhang, P. Finckenberg-Broman, T. Hoang, S. Pan, Z. Xing, M. Staples, and X. Xu, "Right to be forgotten in the era of large language models: Implications, challenges, and solutions," *AI and Ethics*, pp. 1–10, 2024.
- [8] AI@Meta, "Llama 3 model card," 2024. [Online]. Available: [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md)
- [9] Y. Yao, X. Xu, and Y. Liu, "Large language model unlearning," *Advances in Neural Information Processing Systems*, vol. 37, pp. 105 425–105 475, 2024.
- [10] 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.
- [11] R. Zhang, L. Lin, Y. Bai, and S. Mei, "Negative preference optimization: From catastrophic collapse to effective unlearning," *arXiv preprint arXiv:2404.05868*, 2024.
- [12] 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.
- [13] 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.
- [14] 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.
- [15] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, "Direct preference optimization: Your language model is secretly a reward model," *Advances in neural information processing systems*, vol. 36, pp. 53 728–53 741, 2023.
- [16] Z. Wang, M. Shao, R. Karn, L. Mankali, J. Bhandari, R. Karri, O. Sinanoglu, M. Shafique, and J. Knechtel, "Salad: Systematic assessment of machine unlearning on llm-aided hardware design," *arXiv preprint arXiv:2506.02089*, 2025.
- [17] Q. Li, S. Hong, J. Gao, X. Zhang, T. Lan, and W. Cao, "AnalogFed: Federated discovery of analog circuit topologies with generative ai," *arXiv preprint arXiv:2507.15104*, 2025.
- [18] J. Gao, W. Cao, J. Yang, and X. Zhang, "Analoggenie: A generative engine for automatic discovery of analog circuit topologies," *arXiv preprint arXiv:2503.00205*, 2025.
- [19] J. Gao, W. Cao, and X. Zhang, "Analoggenie-lite: Enhancing scalability and precision in circuit topology discovery through lightweight graph modeling," in *Forty-second International Conference on Machine Learning*.
- [20] Z. Dong, W. Cao, M. Zhang, D. Tao, Y. Chen, and X. Zhang, "Cktgnn: Circuit graph neural network for electronic design automation," *arXiv preprint arXiv:2308.16406*, 2023.
- [21] S. Wang, Q. Li, H. He, J. Gao, Z. Wang, Y. Sun, X. Zhang, T. Chi, and W. Cao, "Invited paper: Multi-agent generative synthesis for analog/rf circuit: from scalable topology generation to efficient inverse design," in *2025 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*, 2025, pp. 1–9.
- [22] S. Poddar, C.-T. Ho, Z. Wei, W. Cao, H. Ren, and D. Z. Pan, "Heart: A hierarchical circuit reasoning tree-based agentic framework for ams design optimization," *arXiv preprint arXiv:2511.19669*, 2025.
- [23] 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.
- [24] 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.
- [25] 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*, vol. 43, no. 10, pp. 3184–3197, 2024.
- [26] 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.
- [27] 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.
- [28] 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," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–7.
- [29] 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.
- [30] 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.
- [31] T. Wu, W. Wu, X. Wang, K. Xu, S. Ma, B. Jiang, P. Yang, Z. Xing, Y.-F. Li, and G. Haffari, "Versicode: Towards version-controllable code generation," *arXiv preprint arXiv:2406.07411*, 2024.
- [32] "Ieee standard for systemverilog—unified hardware design, specification, and verification language," *IEEE Std 1800-2023 (Revision of IEEE Std 1800-2017)*, pp. 1–1354, 2024.
- [33] "Ieee standard for verilog hardware description language," *IEEE Std 1364-2005 (Revision of IEEE Std 1364-2001)*, pp. 1–590, 2006.
- [34] S. Palnitkar, *Verilog HDL: a guide to digital design and synthesis*. Prentice Hall Professional, 2003, vol. 1.
- [35] M. Nasr, N. Carlini, J. Hayase, M. Jagielski, A. F. Cooper, D. Ippolito, C. A. Choquette-Choo, E. Wallace, F. Tramèr, and K. Lee, "Scalable extraction of training data from (production) language models," *arXiv preprint arXiv:2311.17035*, 2023.
- [36] 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.
- [37] B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez *et al.*, "Code llama: Open foundation models for code," *arXiv preprint arXiv:2308.12950*, 2023.
- [38] D. Y. Z. X. K. D. W. Z. G. C. X. B. Y. W. Y. L. F. L. Y. X. W. L. Daya Guo, Qihao Zhu, "Deepseek-coder: When the large language model meets programming – the rise of code intelligence," 2024. [Online]. Available: <https://arxiv.org/abs/2401.14196>
- [39] B. Hui, J. Yang, Z. Cui, J. Yang, D. Liu, L. Zhang, T. Liu, J. Zhang, B. Yu, K. Lu *et al.*, "Qwen2. 5-coder technical report," *arXiv preprint arXiv:2409.12186*, 2024.
- [40] Y. Liang, Q. Li, S. Wang, and W. Cao, "When forgetting builds reliability: Llm unlearning for reliable hardware code generation," *arXiv preprint arXiv:2512.05341*, 2025.