VeriReason: Reinforcement Learning with Testbench Feedback for Reasoning-Enhanced Verilog Generation

Yiting Wang^{1,*}, Guoheng Sun^{1,*}, Wanghao Ye¹, Gang Qu¹, Ang Li^{1,†}

¹Department of Electrical Engineering, University of Maryland, Maryland, United States

*Equal contribution

†Corresponding author: angliece@umd.edu

Abstract

Automating Register Transfer Level (RTL) code generation using Large Language Models (LLMs) offers substantial promise for streamlining digital circuit design and reducing human effort. However, current LLM-based approaches for RTL code generation face significant challenges. Methods such as supervised fine-tuning (SFT), in-context learning, and chain-of-thought (CoT) struggle with several critical limitations in the RTL domain: the scarcity of high-quality training data, poor alignment between natural language specifications and generated code, lack of built-in verification mechanisms, and difficulty balancing between model generalization and domain specialization. Inspired by groundbreaking research such as DeepSeek-R1, which combines reinforcement learning with reasoning capabilities, we introduce VeriReason, a comprehensive framework that integrates supervised fine-tuning with Guided Reward Proximal Optimization (GRPO) reinforcement learning specifically tailored for RTL code generation. Using our curated highquality training examples alongside a feedback-driven reward model, VeriReason combines testbench evaluations with structural heuristics to improve specificationcode alignment and eliminate hallucinations. Iterative GRPO embeds intrinsic self-checking and reasoning capabilities, enabling the model to autonomously detect and correct functional errors. On the VerilogEval Benchmark, VeriReason delivers significant improvements: achieving 83.1% functional correctness on the VerilogEval Machine benchmark, substantially outperforming both comparablesized models and much larger commercial systems like GPT-4 Turbo. Additionally, our approach demonstrates up to a 2.8× increase in first-attempt functional correctness compared to baseline methods and exhibits robust generalization to unseen designs. To our knowledge, VeriReason represents the first system to successfully integrate explicit reasoning capabilities with reinforcement learning for Verilog generation, establishing a new state-of-the-art for automated RTL synthesis. The models and datasets are available at:

https://huggingface.co/collections/AI4EDA-CASE Code is Available at: https://github.com/NellyW8/VeriReason

1 Introduction

Register Transfer Level (RTL) code generation is a critical yet labor-intensive task in digital circuit design, directly impacting the efficiency, performance, and power consumption of hardware systems. Traditionally, hardware engineers manually craft RTL code using hardware description languages (HDLs) such as Verilog, which differs significantly from general-purpose programming languages due to its concurrent and structural nature. Recent advancements in large language models (LLMs) offer promising opportunities to automate RTL code generation, substantially reducing the manual

effort and domain expertise required. Leveraging LLMs for RTL generation can accelerate design cycles, minimize human-induced errors, and allow engineers to focus on high-level architectural decisions rather than intricate coding details.

Despite these advantages, LLM-based RTL synthesis encounters three core challenges. First, **data scarcity**: high-quality Verilog examples—and especially paired testbenches or reasoning annotations—are rare, limiting both pretraining and supervised fine-tuning (SFT) and hampering generalization. Second, **weak natural language—code alignment**: LLMs often produce syntactically valid but functionally incorrect Verilog, misinterpreting user specifications and hallucinating invalid structural heuristics (*e.g.*, port matching, net connectivity). Third, **low first-attempt accuracy without self-checking**: current models lack intrinsic mechanisms to detect or correct their own errors, relying instead on external testbench or syntax feedback for iterative refinement. Lastly, **lack of complex logical capability:** Traditional LLMs struggle to handle the intricate interdependencies between components in hardware design, often failing to maintain consistency across module interfaces, state machines, and timing constraints. Without systematic reasoning about component relationships, models produce circuits with logical inconsistencies or incomplete implementations that meet superficial requirements but fail under comprehensive verification.

Recent advances in reasoning and reinforcement learning (RL) have introduced promising approaches to overcome these challenges. Reasoning-augmented models, such as those leveraging chain-of-thought prompting or iterative refinement, have demonstrated the ability to follow multi-step logical patterns, making them particularly suitable for hardware description languages like Verilog that require strict structural correctness and functional dependencies. These reasoning mechanisms help LLMs better understand circuit intent and adhere to design constraints, and can better ensure the alignment between natural language and result. Methods such as Guided Reward Proximal Optimization (GRPO)[16] combine the strengths of SFT with reward-driven RL, enabling models to learn effectively even with minimal data and explicit feedback. By employing RL-based strategies, LLMs are trained not merely on predicting the next token but on achieving specific, meaningful outcomes, thus improving their logical reasoning, alignment, and self-checking capabilities.

Our Proposed Framework. To address the challenges in RTL generation with LLMs, we propose a novel framework, VeriReason, combining supervised fine-tuning (SFT) and GRPO reinforcement learning, specifically tailored for Verilog RTL generation with a specially designed dataset featuring reasoning steps and testbenches. Our approach systematically tackles four critical limitations that hinder existing LLM-based hardware design methods: data scarcity in domain-specific code, natural language-code alignment issues, lack of self-checking behavior, and insufficient complex logical capabilities. Each of these challenges requires specialized techniques that we incorporate into the VeriReason framework, as detailed in the following sections.

Data Scarcity in Domain-Specific Code: We introduce a reasoning-distillation and testbench-generation pipeline to augment existing prompt—code pairs with high-quality testbenches and human-style reasoning steps, producing a high-quality dataset. Furthermore, we demonstrate that even with as few as 20 annotated examples from the VeriReason dataset, GRPO yields substantial performance gains, dramatically lowering the bar for required training data.

Natural Language-Code Alignment: VeriReason employ a reward model that evaluates generated Verilog code against specifications using feedback from structural heuristics. Through GRPO optimization, the model learns to internalize structural constraints, effectively reducing hallucinations and ensuring structural correctness by penalizing invalid constructs across both interface definitions and internal hierarchy of the circuit design.

Lack of Self-Checking Behavior: Our reinforcement learning framework inherently encourages the model to develop self-checking capabilities by iteratively refining outputs based on testbench feedback-driven rewards. Over training iterations, the model learns to anticipate and rectify errors internally, significantly enhancing first-attempt functional correctness.

Lack of complex logical capabilities: VeriReason incorporates explicit reasoning steps throughout the design process, requiring the model to articulate its design decisions and verify logical consistency before implementation. By decomposing complex circuit specifications into manageable conceptual components and reasoning about their interactions, the model develops more coherent and complete implementations that maintain logical integrity across the entire design.

Our key contributions are summarized as follows:

- We design a novel framework VeriReason, which integrates supervised fine-tuning with GRPO-based reinforcement learning and reasoning-augmented design processes for Verilog RTL code generation.
- Our approach addresses critical shortcomings in RTL generation through a reasoningdistillation pipeline for data scarcity, reward-driven structural evaluation for NL-code alignment, testbench feedback mechanisms for self-checking behavior, and explicit reasoning steps for handling complex logical dependencies.
- We create a high quality dataset with reasoning and testbench that would be open-sourced to the benefit community.
- The framework achieves state-of-the-art performance in RTL generation tasks, demonstrating substantial improvements in first-attempt functional correctness, structural validity, and generalization capabilities with minimal training data. It delivers up to a 2.8× increase in first-attempt functional correctness compared to baseline models while outperforming existing state-of-the-art methods across multiple benchmarks. This improvement is particularly notable in smaller parameter models. Remarkably, the framework achieves 83.1% pass@5 on VerilogEval-Machine, even surpassing much larger models including GPT-4 Turbo.

2 Background

Recent years have seen a surge of interest in applying large language models (LLMs) to hardware design, particularly for generating Register-Transfer Level (RTL) code in Verilog.[7, 9, 2, 1, 20, 17, 25, 11, 10, 18, 13, 4]. Previous research have shown significant potential for LLM-based RTL generation in automating parts of the hardware design process using finetuning, or using differnt prompting techniques. However, it has also revealed several key challenges of produce correct and efficient hardware designs reliably.

Challenges in LLM-based RTL generation tasks Researchers investigated fine-tuning LLMs using domain-specific data and techniques to improve their performance [9, 14]. For example, Liu et.al introduced ChipNeMo [7], which fine-tunes a general-purpose LLM on internal NVIDIA datasets for various chip design tasks. Similarly, Thankur et.al developed VeriGen [17] to improve Verilog generation capabilities. Subsequent works, such as RTLCoder [9] is trained based on automatically generated datasets. BetterV [14], finetunes LLM by converting Verilog code to the C language. While effective, these methods face challenges of scalability and generalizability due to their high demand for high quality instruction-code pairs. The inherent limitation of lack of real RTL code and the low quality of generated code make it hard to make further improvement.

Moreover, LLM generated Verilog code often face the issue of hallucination. Prompt-based methods [1, 2] rely heavily on the quality and clarity of the input prompts, facing difficulties in consistently aligning complex, multi-step circuit specifications with the generated code. These methods often suffer from hallucinations or syntactically correct yet functionally incorrect outputs due to inadequate contextual understanding. Moreover, they inherently lack iterative refinement capabilities, making them incapable of progressively improving RTL code quality. Chain of Thought (CoT) methods [24], which encourage models to generate step-by-step reasoning sequences, typically excel in structured reasoning tasks but face challenges in RTL contexts due to the strict requirement for functional correctness and structural precision. Although CoT enhances reasoning, its effectiveness heavily depends on the clarity and correctness of intermediate reasoning steps, which can still suffer from errors in the absence of explicit correctness feedback mechanisms. Furthermore, the CoT has made the process of generation very ineffective due to long inference time.

2.1 Reinforcement Learning for LLM Reasoning

Recent research has demonstrated the potential of reinforcement learning (RL) techniques to significantly enhance the reasoning capabilities of large language models (LLMs) [12, 16, 5]. By providing explicit rewards for logical correctness and step-wise reasoning, RL enables models to autonomously discover effective problem-solving strategies, often mirroring structured human reasoning [21, 23]. Applications span mathematical problem solving (where RL fine-tuning on step-by-step correctness or final answer accuracy yields substantial improvements [16, 5]) and code generation, where preference optimization and RL from feedback have led to greater code validity and efficiency [3].

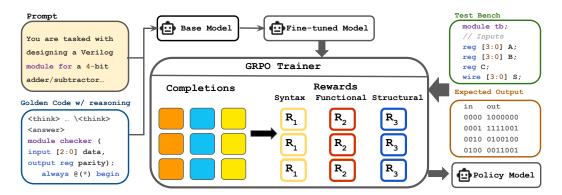


Figure 1: Workflow of VeriReason. The framework combines supervised fine-tuning with GRPO reinforcement learning. A base model is fine-tuned and then improved through the GRPO trainer, which leverages multiple reward signals (syntax, functional, and structural correctness) derived from testbench execution and code analysis. The model incorporates explicit reasoning (<think> blocks) to break down complex hardware design tasks.

Most successful approaches build upon policy gradient algorithms such as Proximal Policy Optimization (PPO) [15] or, more recently, Group Relative Policy Optimization (GRPO) [16, 6]. GRPO, in particular, compares groups of generated responses rather than evaluating them in isolation, enabling the model to build a deeper understanding of what constitutes high-quality reasoning through relative comparisons. The effectiveness of these frameworks depends on carefully designed reward functions that accurately reflect the target domain. For hardware description tasks like Verilog generation, structural similarity as provides a clear, unambiguous reward signal [19], encouraging models to internalize domain-specific constraints and develop robust reasoning capabilities that translate natural language specifications into golden-code similarity.

3 Methodology

We present VeriReason, a comprehensive framework that combines supervised fine-tuning with reinforcement learning specifically tailored for Verilog RTL generation. Our approach addresses four critical challenges in automated hardware design: (1) data scarcity of high-quality RTL examples, (2) weak alignment between natural language specifications and generated code, (3) lack of self-checking mechanisms in current models, and (4) insufficient complex logical reasoning capabilities for hardware design.

As shown in Figure 1, VeriReason employs a multi-stage approach to generate high-quality Verilog code. First, a base model is fine-tuned on a curated dataset consisting of high-quality prompt-code pairs enhanced with explicit reasoning steps and testbenches. This fine-tuned model produces initial code implementations that are then evaluated through our reward system, which combines three key components: syntax correctness, functional validation via testbench execution, and structural analysis. The GRPO (Guided Reward Proximal Optimization) trainer leverages these rewards to iteratively improve the model's ability to generate correct code while maintaining alignment with the original specification. Through this process, the model learns to incorporate reasoning steps (shown as <think> blocks) that decompose complex hardware design problems into manageable components, resulting in a policy model capable of generating functionally correct and structurally sound Verilog implementations on the first attempt.

3.1 Reinforcement Learning Framework

Our approach adapts Guided Reward Proximal Optimization (GRPO) specifically for RTL code generation. Unlike traditional RL methods, our framework incorporates domain-specific constraints and verification mechanisms directly into the learning process, providing immediate feedback on functional correctness, syntax, and specification adherence. This targeted optimization enables the model to efficiently learn correct Verilog implementation patterns while minimizing hallucinations and specification misalignments.

3.1.1 Group Relative Policy Optimization (GRPO)

We adopt Group Relative Policy Optimization (GRPO) as our core reinforcement learning algorithm due to its efficiency and demonstrated effectiveness in tasks requiring complex reasoning. GRPO provides several advantages over traditional reinforcement learning methods like Proximal Policy Optimization (PPO), including lower memory requirements and more stable training dynamics.

In GRPO, the language model serves as the policy network, taking a natural language specification q as input and producing a sequence of tokens representing Verilog code as actions. The policy distribution factors across tokens: $\pi_{\theta}(a|q) = \prod_{t=1}^{N} \pi_{\theta}(a_t|q, a_{< t})$, where π_{θ} represents the policy parameterized by θ , a is the complete sequence of tokens (the Verilog code), and a_t is the token at position t.

Unlike PPO, which requires a separate value function, GRPO estimates advantages using group-based sampling. For each natural language specification q, we generate a group of G candidate Verilog implementations $\{o_1, o_2, \ldots, o_G\}$ from the current policy and compute rewards for each. The GRPO objective function is defined as:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)} \left[\frac{1}{G} \sum_{i=1}^G \min \left(r_i \cdot \rho_i, \text{clip}(\rho_i, 1 - \epsilon, 1 + \epsilon) \cdot r_i \right) \right] - \beta \cdot D_{\text{KL}}(\pi_{\theta}(\cdot | q) \| \pi_{\text{ref}}(\cdot | q))$$
(1)

where: where $\rho_i = \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{\text{old}}}(o_i|q)}$ is the importance sampling ratio, r_i is the normalized reward for candidate o_i , ϵ is a hyperparameter controlling the clipping range, β is a coefficient balancing the KL divergence penalty, π_{ref} is a reference policy (typically the supervised fine-tuned model), and D_{KL} is the Kullback-Leibler divergence.

The advantage estimation in GRPO is simplified by normalizing rewards within each group, where $r_i = \frac{R(o_i) - \mu_R}{\sigma_R + \delta}$, with $R(o_i)$ being the raw reward for output o_i , μ_R and σ_R are the mean and standard deviation of rewards within the group, and δ is a small constant for numerical stability.

This group-based normalization provides several benefits: it eliminates the need for a separate value network, reduces variance in advantage estimation, and naturally compares alternative implementations of the same specification, which aligns well with the goal of generating functionally correct Verilog code.

3.2 Reward Model

Our reward function combines both structural correctness and functional validation to provide comprehensive feedback during training. The reward R for a generated Verilog implementation is computed as:

$$R(o) = \begin{cases} 2.0, & \text{if functionally correct} \\ 0.1 + 1.0 \cdot \text{AST}_{\text{score}}(o), & \text{if syntactically correct} \\ 0, & \text{otherwise} \end{cases}$$
 (2)

where: Functional correctness is determined by running the generated code through testbenches and comparing outputs with the expected behavior. Syntactic correctness is verified by successful parsing of the Verilog code. The $AST_{score}(o)$ measures structural similarity between the generated code's Abstract Syntax Tree (AST) and reference implementations, with values ranging from 0 to 1.

The AST score provides a fine-grained measure of structural correctness even when the code is not functionally perfect. VeriReason employs a hierarchical AST comparison algorithm specifically tailored for Verilog code structures, where ASTscore(o) = $\sum_{c \in C} w_c \cdot (0.6 \cdot \sin_c + 0.5 \cdot \cos_c - 0.3 \cdot \text{red}_c)$ calculates weighted structural similarity across categories $C = \{\text{module, port, always, ...}\}$ with respective importance weights w_c . For each category c, we compute sequence similarity \sin_c using Levenshtein distance, coverage $\cos_c = |G_c \cap D_c|/|G_c|$ between generated elements C_c and golden elements C_c , and redundancy $\cos_c = |D_c - G_c|/|D_c|$ to penalize hallucinated structures. This domain-specific structural analysis enables our model to maintain correct interface definitions,

signal declarations, and control logic while internalizing hardware design patterns during GRPO optimization.

For functional verification, we use testbenches to evaluate the generated Verilog against reference implementations. A generated design is considered functionally correct only when it passes all test cases in the testbench, providing identical output signals to those of the golden reference for all test vectors.

3.3 Data Preprocessing

We address the critical issue of data scarcity in RTL generation, and the challenge of low dataset quality through a data augmentation pipelines, and a data filtration pipeline.

3.3.1 Data Filtration

We implement a two-stage adaptive filtration process to optimize the dataset for GRPO training effectiveness. For each sample s in our initial dataset D, we generate a set of k=8 candidate implementations $\{o_1,o_2,\ldots,o_k\}$ using our base model and compute their corresponding rewards $\{r_1,r_2,\ldots,r_k\}$ based on our reward function R defined in Equation 2. The filtration process is formalized as follows: $D_{\text{filtered}}=\{s\in D\mid \mu_r(s)\in [\alpha_{\min},\alpha_{\max}] \text{ and } \sigma_r(s)>\beta\}$, where: $\mu_r(s)=\frac{1}{k}\sum_{i=1}^k r_i$ is the mean reward for sample s, $\sigma_r(s)=\sqrt{\frac{1}{k}\sum_{i=1}^k (r_i-\mu_r(s))^2}$ is the standard deviation of rewards, $\alpha_{\min}=0.3$ is the minimum acceptable mean reward, $\alpha_{\max}=1.8$ is the maximum acceptable mean reward, $\beta=0.1$ is the minimum acceptable reward variance. This filtration strategy excludes samples that are either too difficult (consistently low rewards) or too trivial (consistently high rewards), retaining only those samples that provide meaningful learning signals for the GRPO algorithm. Specifically, we filter out data where all rewards are zero or the average reward is below α_{\min} , as well as samples where all generations achieve near-perfect scores.

Additionally, we compute a difficulty score $\delta(s)$ for each remaining sample:

$$\delta(s) = 1 - \frac{\mu_r(s) - \alpha_{\min}}{\alpha_{\max} - \alpha_{\min}}$$
(3)

Samples are then categorized into "simple" and "hard" portions based on this difficulty score, with samples where $\delta(s)>0.5$ classified as "hard" and the remainder as "simple." This stratification enables targeted training strategies that progressively build model competence. The final dataset comprises 1149 samples in the hard level and 743 samples in the easy level.

3.3.2 Reasoning Generation with Optimization

To enhance the model's reasoning capabilities, we augment the training data with explicit reasoning steps that decompose complex hardware design problems into manageable components. Our reasoning generation pipeline extracts natural language specifications and corresponding Verilog implementations from the original dataset and uses chain-of-thought prompting with domain-specific guidance to generate detailed reasoning steps that explain the design choices, module interfaces, and implementation details. We apply an optimization process where the model is encouraged to critique its own reasoning and suggest improvements. When the model identifies a better alternative solution, we generate an improved implementation and incorporate it into the dataset. This self-improvement mechanism enables the model to iteratively refine both its reasoning process and the quality of generated code. The resulting dataset contains not only input-output pairs but also explicit reasoning traces that guide the model to develop stronger internal reasoning capabilities.

3.3.3 Testbench Generation

To provide reliable functional correctness signals during training, we develop an automated testbench generation pipeline. The pipeline analyzes the input-output specifications to identify signal characteristics, boundary conditions, and expected behaviors. It generates comprehensive testbenches covering both typical and edge cases, applying multiple test vector generation strategies, including directed testing for explicit requirements and constrained random testing for broader coverage. For each specification, our system generates at least 100 diverse test vectors to ensure adequate functional

coverage and validates testbenches by confirming they correctly identify known-good and known-bad implementations. The generated testbenches are used both for reward computation during reinforcement learning and for final validation of model outputs. This approach ensures that the model learns to produce not just syntactically correct but functionally valid Verilog implementations aligned with the original specifications.

The combination of these data augmentation techniques with GRPO-based reinforcement learning creates a powerful framework for RTL generation that addresses the challenges of data scarcity, weak natural language-code alignment, and lack of self-checking behaviors. By encoding hardware design best practices through both structural and functional rewards, VeriReason enables the model to internalize domain-specific constraints and develop robust reasoning capabilities for RTL synthesis.

4 Experiments

4.1 Experimental Setup

Our primary dataset is derived from the RTLCoder [9] dataset of 26500 samples. We apply a thorough filtration technique on the dataset. First, we apply a simple syntax check to ensure we keep only the syntactically valid code. Then we use ChatGPT-4.1 to check whether the code matches the input prompt correctly, and generate reasoning steps for the code. For code that does not fully match the input prompt, it is re-generated and checked for syntax. We then generate a testbench for each entry with at least 100 test cases for best coverage.

The testbench output is also saved to the dataset; in this way, during the GRPO, the testbench will only run once on the generated code, and the outcome will be directly compared to the golden code's output. Next, we run the dataset on the Qwen2.5B model to generate the code and its corresponding rewards. Based on the reward, we split the dataset into simple and hard portions for the next training steps. We end up with 1149 samples in the hard level and 743 samples in the easy level.

Our evaluation focuses on the primary benchmarks for Verilog code generation, VerilogEval [8]. The comprehensive benchmark containing both machine-generated and human-crafted Verilog specifications. VerilogEval-Machine contains 143 samples with algorithmically generated specifications, while VerilogEval-Human includes 156 samples with human-written specifications.

We evaluate VeriReason across multiple model scales to assess parameter efficiency, including Qwen2.5-1.5B, Qwen2.5-3B, Qwen2.5-7B, and CodeLlama3-7B architectures. GRPO is used as our default RL algorithm. RL-specific settings include a generation temperature of 0.5, a total batch size of 16 (8 rollouts each), an update batch size of 2 per GRPO step with a gradient accumulation of 8, a lowered learning rate of 1.0e-6 with constant scheduler type, and repetition penalty of 1.3. The reward model follows the design in Equation 2, with execution correctness verified using industry-standard Verilog simulators, Iverilog [22].

4.2 Main Results

Table 1 compares VeriReason against state-of-the-art Verilog generation models. Our approach achieves superior performance across all model sizes, with VeriReason-Qwen2.5-7B demonstrating remarkable gains over base models: +17.1 and +24.0 percentage points on pass@1 for VerilogEval-Machine and VerilogEval-Human, respectively.

Even our smallest model, VeriReason-Qwen2.5-1.5B, outperforms many larger models after applying our reinforcement learning framework. The performance gains are particularly notable for smaller models, demonstrating the effectiveness of our approach in optimizing parameter efficiency. We observe that VeriReason-Qwen2.5-7B achieves state-of-the-art performance on Machine pass@5 (83.1%), outperforming even GPT-4-Turbo on this metric, despite having significantly less parameters.

Furthermore, all VeriReason models establish themselves as the **best performers** in their respective parameter size categories (1.5B, 3B, and 7B), highlighting the robustness of our approach across model scales. The substantial improvements observed in VeriReason-codeLlama-7B (+25.2 percentage points in Machine pass@1) further demonstrates that our method generalizes effectively across different model architectures.

Table 1: Comparative analysis of Verilog code generation performance. Gray highlighting denotes the overall state-of-the-art results. (**bold**) indicates the best results in the model size category. Color-coded numbers show performance deltas relative to base models (green: improvement).

Category	Method	Params.	Open Source	VerilogEval-Machine		VerilogEval-Human	
				pass@1	pass@5	pass@1	pass@5
Base Model	GPT-3.5-Turbo	N/A	Х	63.5	78.0	31.2	47.0
	GPT-4o-mini	N/A	X	66.0	72.4	54.2	62.0
	GPT-4-Turbo	N/A	×	72.5	<u>83.0</u>	64.3	76.1
	Qwen-2.5-Coder	1.5B	✓	25.6	40.8	8.3	17.9
	Qwen2.5-Coder	3B	✓	48.4	58.9	21.3	32.7
	Qwen2.5-Coder	7B	✓	52.7	69.7	23.9	41.1
	CodeLlama	7B	✓	26.1	49.1	18.8	28.6
	CodeQwen1.5-7B-Chat	7B	✓	29.1	61.9	14.8	36.8
	DeepSeek-Coder	6.7B	\checkmark	8.8	34.3	4.9	19.3
	DeepSeek-V3	671B	\checkmark	79.2	80.7	66.1	72.1
Fine-tuned Generation	ChipNeMo†	70B	Х	53.8	N/A	27.6	N/A
	BetterV-CodeOwen†	7B	Х	68.1	79.4	46.1	53.7
	RTLLLM [†]	13B	X	65.3	77.2	43.7	51.8
	VerilogEval [†]	16B	X	46.2	67.3	28.8	45.9
	VeriGen [†]	16B	✓	44.0	52.6	30.3	43.9
	RTLCoder-DeepSeek-Coder	6.7B	✓	$37.2_{+28.4}$	$64.9_{+30.6}$	$16.9_{+12.0}$	$35.7_{+16.4}$
VeriReason (Ours)	VeriReason-Qwen2.5-1.5B	1.5B	✓	44.7 _{+19.1}	49.1 _{+8.3}	23.5 _{+15.2}	26.7+8.8
	VeriReason-Qwen2.5-3B	3B	✓	55.9 _{+7.5}	72.8 _{+13.9}	33.2 _{+11.9}	47.4 _{+14.7}
	VeriReason-Qwen2.5-7B	7B	✓	69.8 _{+17.1}	83.1+13.4	47.9 _{+24.0}	58.4 _{+17.3}
	VeriReason-codeLlama-7B	7B	\checkmark	51.3+25.2	64.0+14.9	$27.5_{+8.7}$	39.9+11.3

^{†:} Reported Results.

4.3 Training Dynamics Analysis

Figure 2 illustrates the training dynamics of our VeriReason models across three different parameter scales (1.5B, 3B, and 7B). We track both the mean reward values (top row) and their standard deviations (bottom row) throughout the reinforcement learning process.

4.3.1 Reward Progression

The reward curves reveal distinct learning patterns across model sizes. The 1.5B model demonstrates steady, monotonic improvement in reward values from approximately 0.5 to 0.8 over 800 training steps, suggesting a consistent optimization path. In contrast, the 3B model exhibits higher variance in its learning trajectory, with reward values fluctuating between 0.6 and 0.8, before ultimately converging to above 0.8 by step 400. The 7B model shows the most pronounced oscillatory behavior, with rewards ranging between 0.50 and 0.65, reflecting the increased complexity of optimizing larger parameter spaces.

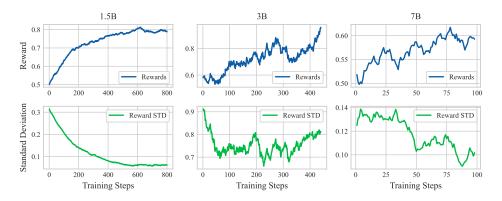


Figure 2: Reward line and std line of VeriReason

4.3.2 Reward Stability

The standard deviation plots (bottom row) provide critical insights into training stability. The 1.5B model demonstrates exceptional convergence properties, with reward variability consistently decreasing from 0.3 to below 0.1 throughout training. This smooth reduction in standard deviation correlates with the steady increase in mean rewards, indicating robust learning. The 3B model presents a more complex pattern, with initial rapid variability reduction followed by fluctuations between 0.7 and 0.8, suggesting periodic exploration-exploitation transitions. The 7B model's standard deviation exhibits the most dynamic behavior, oscillating between 0.09 and 0.14, which aligns with its more variable reward progression.

Interestingly, despite having fewer parameters, the 1.5B model achieves the most stable convergence pattern, with monotonically decreasing standard deviation. This suggests that smaller models may benefit more consistently from our reinforcement learning framework, while larger models engage in more extensive exploration of the parameter space before convergence. The final standard deviation values (approximately 0.05 for 1.5B, 0.8 for 3B, and 0.1 for 7B) indicate that all models eventually reach stable policy configurations, though with different convergence trajectories.

This analysis provides empirical evidence that our reinforcement learning approach effectively optimizes models across different parameter scales, with larger models requiring more complex optimization paths but ultimately achieving higher reward values, consistent with their superior performance on the VerilogEval benchmarks shown in Table 1.

4.3.3 The effect of SFT and GRPO

Table 2 presents a systematic analysis of how each training component—Supervised Fine-Tuning (SFT) and Group Relative Policy Optimization (GRPO)—contributes to VeriReason's performance across different model sizes. The results demonstrate that SFT provides substantial initial performance gains across all model architectures, with smaller models showing the most dramatic relative improvements. This suggests that SFT effectively addresses the challenge of limited domain knowledge by providing high-quality training examples with explicit reasoning steps. The addition of GRPO further enhances performance across all models and

Table 2: Ablation studies results.

Model	Training Stage	VerilogEv	al-Machine	VerilogEval-Human		
Model		pass@1	pass@5	pass@1	pass@5	
Qwen2.5-1.5B	Base	25.6	40.8	8.2	17.9	
	+ SFT	38.6	46.3	17.8	23.9	
	+ GRPO	44.7	49.1	23.5	26.7	
Qwen2.5-3B	Base	48.4	58.9	21.3	32.7	
	+ SFT	51.9	69.9	31.3	45.1	
	+ GRPO	55.9	72.8	33.2	47.4	
Qwen2.5-7B	Base	52.7	69.7	23.9	41.1	
	+ SFT	63.4	79.9	43.4	56.2	
	+ GRPO	69.8	83.1	47.9	58.4	
CodeLlama-7B	Base	26.1	49.1	18.8	28.6	
	+ SFT	41.1	58.3	23.2	31.5	
	+ GRPO	51.3	64.0	27.5	39.9	

benchmarks, demonstrating the value of reinforcement learning with testbench feedback. This synergistic relationship shows that while SFT provides the foundation for understanding Verilog syntax and semantics, GRPO enhances the model's ability to produce functionally correct implementations by internalizing structural constraints and developing self-checking capabilities.

5 Conclusion

This paper presents VeriReason, a comprehensive framework integrating supervised fine-tuning with GRPO-based reinforcement learning for Verilog RTL generation. By combining explicit reasoning with testbench-driven feedback, our approach addresses key challenges in LLM-based hardware design: data scarcity, weak language-code alignment, lack of self-checking behavior, and insufficient logical reasoning. While VeriReason achieves state-of-the-art performance across model scales, the approach incurs significant computational overhead during both training (requiring numerous testbench evaluations per iteration) and inference (where reasoning steps increase generation time by 2.5-3×). Despite these limitations, the consistent improvements across different architectures demonstrate our method's robustness and transferability. Our reward model's integration of structural correctness and functional validation encourages models to develop intrinsic self-checking capabilities—a crucial advancement for autonomous hardware design. By open-sourcing our models and datasets, we aim to accelerate progress in LLM-based hardware design, transform digital circuit development practices, and inspire future work on computational efficiency.

References

- [1] Jason Blocklove, Siddharth Garg, Ramesh Karri, and Hammond Pearce. Chip-chat: Challenges and opportunities in conversational hardware design. In 5th ACM/IEEE Workshop on Machine Learning for CAD, MLCAD. IEEE, 2023.
- [2] Kaiyan Chang, Ying Wang, Haimeng Ren, Mengdi Wang, Shengwen Liang, Yinhe Han, Huawei Li, and Xiaowei Li. Chipgpt: How far are we from natural language hardware design. CoRR, abs/2305.14019, 2023.
- [3] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv*, 2107.03374, July 2021.
- [4] Yonggan Fu, Yongan Zhang, Zhongzhi Yu, Sixu Li, Zhifan Ye, Chaojian Li, Cheng Wan, and Yingyan Celine Lin. GPT4AIGChip: Towards Next-Generation AI Accelerator Design Automation via Large Language Models. arXiv, 2309.10730, September 2023. Accepted by ICCAD 2023.
- [5] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv*, 2501.12948, January 2025.
- [6] Nathan Lambert, Salman Khan, Sunabha Chatterjee, Shandong Wu, William Mitchell, Yuntian Deng, Hang Gao, Saurav Pahadia, Roshni Sahoo, Xuechen Li, Anca Dragan, and Jacob Steinhardt. Reinforcement learning with verifiable rewards: Grpo's effective loss, dynamics, and success amplification. *arXiv*, 2503.06639, March 2024.
- [7] Mingjie Liu, Teodor-Dumitru Ene, Robert Kirby, Chris Cheng, Nathaniel Ross Pinckney, Rongjian Liang, Jonah Alben, Himyanshu Anand, Sanmitra Banerjee, Ismet Bayraktaroglu, Bonita Bhaskaran, Bryan Catanzaro, Arjun Chaudhuri, Sharon Clay, Bill Dally, Laura Dang, Parikshit Deshpande, Siddhanth Dhodhi, Sameer Halepete, Eric Hill, Jiashang Hu, Sumit Jain, Brucek Khailany, Kishor Kunal, Xiaowei Li, Hao Liu, Stuart F. Oberman, Sujeet Omar, Sreedhar Pratty, Jonathan Raiman, Ambar Sarkar, Zhengjiang Shao, Hanfei Sun, Pratik P. Suthar, Varun Tej, Kaizhe Xu, and Haoxing Ren. Chipnemo: Domain-adapted Ilms for chip design. *CoRR*, abs/2311.00176, 2023.
- [8] Mingjie Liu, Nathaniel Ross Pinckney, Brucek Khailany, and Haoxing Ren. Invited paper: Verilogeval: Evaluating large language models for verilog code generation. In *IEEE/ACM International Conference on Computer Aided Design, ICCAD*. IEEE, 2023.
- [9] Shang Liu, Wenji Fang, Yao Lu, Jing Wang, Qijun Zhang, Hongce Zhang, and Zhiyao 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.
- [10] Shang Liu, Wenji Fang, Yao Lu, Jing Wang, Qijun Zhang, Hongce Zhang, and Zhiyao Xie. Rtlcoder: Fully open-source and efficient llm-assisted rtl code generation technique, 2024.
- [11] Yao Lu, Shang Liu, Qijun Zhang, and Zhiyao Xie. Rtllm: An open-source benchmark for design rtl generation with large language model. In *Proceedings of the 29th Asia and South Pacific Design Automation Conference*, ASPDAC '24. IEEE Press, 2024.
- [12] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [13] Hammond Pearce, Benjamin Tan, and Ramesh Karri. Dave: Deriving automatically verilog from english. In *Proceedings of the 2020 ACM/IEEE Workshop on Machine Learning for CAD*, pages 27–32, 2020.
- [14] Zehua Pei, Hui-Ling Zhen, Mingxuan Yuan, Yu Huang, and Bei Yu. Betterv: Controlled verilog generation with discriminative guidance. In *Forty-first International Conference on Machine Learning, ICML*. OpenReview.net, 2024.

- [15] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv*, 1707.06347, July 2017.
- [16] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv*, 2402.03300, February 2024.
- [17] Shailja Thakur, Baleegh Ahmad, Hammond Pearce, Benjamin Tan, Brendan Dolan-Gavitt, Ramesh Karri, and Siddharth Garg. Verigen: A large language model for verilog code generation. *ACM Trans. Design Autom. Electr. Syst.*, 2024.
- [18] Yunda Tsai, Mingjie Liu, and Haoxing Ren. Rtlfixer: Automatically fixing RTL syntax errors with large language model. In *Proceedings of the 61st ACM/IEEE Design Automation Conference*, DAC. ACM, 2024.
- [19] Ning Wang, Bingkun Yao, Jie Zhou, Xi Wang, Zhe Jiang, and Nan Guan. Large language model for verilog generation with code-structure-guided reinforcement learning, 2025.
- [20] Yiting Wang, Wanghao Ye, Ping Guo, Yexiao He, Ziyao Wang, Yexiao He, Bowei Tian, Shwai He, Guoheng Sun, Zheyu Shen, Sihan Chen, Ankur Srivastava, Qingfu Zhang, Gang Qu, and Ang Li. Symrtlo: Enhancing rtl code optimization with llms and neuron-inspired symbolic reasoning, 2025.
- [21] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *arXiv*, 2201.11903, January 2022.
- [22] Stephen Williams. The icarus verilog compilation system. [Online], 2023. Available: https://github.com/steveicarus/iverilog.
- [23] Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. Logic-rl: Unleashing llm reasoning with rule-based reinforcement learning. *arXiv*, 2504.04253, April 2025.
- [24] Yiyao Yang, Fu Teng, Pengju Liu, Mengnan Qi, Chenyang Lv, Ji Li, Xuhong Zhang, and Zhezhi He. Haven: Hallucination-mitigated llm for verilog code generation aligned with hdl engineers, 2025.
- [25] Ruizhe Zhong, Xingbo Du, Shixiong Kai, Zhentao Tang, Siyuan Xu, Hui-Ling Zhen, Jianye Hao, Qiang Xu, Mingxuan Yuan, and Junchi Yan. Llm4eda: Emerging progress in large language models for electronic design automation. *arXiv preprint arXiv:2401.12224*, 2023.