# VeriMoA:
# A Mixture-of-Agents Framework for Spec-to-HDL Generation

**Heng Ping** [1]  **Arijit Bhattacharjee** [2]  **Peiyu Zhang** [1]  **Shixuan Li** [1]  **Wei Yang** [1]  **Anzhe Cheng** [1]  **Xiaole Zhang** [1]
**Jesse Thomason** [1]  **Ali Jannesari** [2]  **Nesreen Ahmed** [3]  **Paul Bogdan** [1]

## ABSTRACT

Automation of Register Transfer Level (RTL) design can help developers meet increasing computational demands. Large Language Models (LLMs) show promise for Hardware Description Language (HDL) generation, but face challenges due to limited parametric knowledge and domain-specific constraints. While prompt engineering and fine-tuning have limitations in knowledge coverage and training costs, multi-agent architectures offer a training-free paradigm to enhance reasoning through collaborative generation. However, current multi-agent approaches suffer from two critical deficiencies: susceptibility to noise propagation and constrained reasoning space exploration. We propose VeriMoA, a training-free mixture-of-agents (MoA) framework with two synergistic innovations. First, a **quality-guided caching mechanism** to maintain all intermediate HDL outputs and enables quality-based ranking and selection across the entire generation process, encouraging knowledge accumulation over layers of reasoning. Second, a **multi-path generation strategy** that leverages C++ and Python as intermediate representations, decomposing specification-to-HDL translation into two-stage processes that exploit LLM fluency in high-resource languages while promoting solution diversity. Comprehensive experiments on VerilogEval 2.0 and RTLLM 2.0 benchmarks demonstrate that VeriMoA achieves 15–30% improvements in Pass@1 across diverse LLM backbones, especially enabling smaller models to match larger models and fine-tuned alternatives without requiring costly training.

## 1  INTRODUCTION

As the semiconductor industry faces unprecedented pressure to deliver high-performance hardware, the automation of Register Transfer Level (RTL) design, a critical component of the chip development pipeline, has become increasingly vital (Liu et al., 2023). The ability to automatically generate accurate Hardware Description Language (HDL) from natural language specifications would dramatically accelerate design cycles and reduce engineering pressure (Lu et al., 2024). Recent advances in Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse domains, including Electronic Design Automation (EDA), where they have shown promising potential for automating HDL generation from natural language instructions (Thakur et al., 2023). Unfortunately, despite their impressive performance in general-purpose programming tasks, applying LLMs to HDL generation presents unique challenges.

Unlike widely-used languages such as C++ or Python, HDL represents a specialized domain with limited representa-

tion in LLM pretraining corpora, resulting in comparatively sparse parametric knowledge for hardware design (Liu et al., 2024b). Moreover, HDL exhibits distinct operational semantics and constraints. For example, HDL must capture concurrent hardware behavior, satisfy timing constraints, and adhere to synthesis requirements, requiring reasoning over temporal and physical limitations inherent to circuit design. These domain-specific characteristics make it difficult for general-purpose LLMs to generate correct and efficient HDL code from specifications (Alsaqer et al., 2025).

To address these limitations, initial research focused on *model-centric* strategies to enhance LLM's HDL generation capabilities. The first wave employed prompt engineering to elicit the LLM's sparse HDL knowledge: ParaHDL (Sun et al., 2025) proposed task-specific prompt paradigms embedding HDL syntax patterns, while AoT (DeLorenzo et al., 2025) used tailored Chain-of-Thought templates for different circuit categories. However, these methods are fundamentally limited by the LLM's inadequate pre-existing HDL understanding. A second wave progressed towards augmenting parametric knowledge through fine-tuning on high-quality HDL datasets. Methods like RTLCoder (Liu et al., 2024a) and AutoVCoder (Gao et al., 2024) curated large-scale, simulation-verified HDL corpora for supervised fine-

---
[1]University of Southern California, Los Angeles, California, USA [2]Iowa State University, Ames, Iowa, USA [3]Cisco AI Research, San Jose, California, USA
Correspondence to: Heng Ping <hping@usc.edu>, Paul Bogdan <pbogdan@usc.edu>

tuning, while VeriSeek (Wang et al., 2024b) and ChipSeek-R1 (Chen et al., 2025) employed reinforcement learning to further enhance reasoning. Despite substantial improvements, the model-centric philosophy faces critical limitations: prompt engineering remains bounded by limited pre-existing knowledge; fine-tuning requires significant data curation and training costs; most fundamentally, both rely on monolithic generation, lacking mechanisms for systematic collaboration and progressive improvement.

Recognizing the risks of a monolithic generation process, an alternative research direction constructs system-level frameworks, with multi-agent architectures emerging as a prominent approach for HDL generation. However, the prevailing collaborative paradigms are fundamentally flawed, ranging from the rigid, linear pipelines of MAGE (Zhao et al., 2024) that propagate errors, to the unstructured debates of Coopetitive V (Mi et al., 2024) that lead to chaotic exploration. This progression exposes two profound deficiencies in current approaches. First, the collaborative processes of these frameworks are highly susceptible to noise and cannot guarantee information gain. Without a principled mechanism to filter out flawed information, the system's knowledge base can be easily corrupted, preventing it from productively building upon prior successes. Second, even when correct information is preserved, they operate within a constrained reasoning space. Both brittle and chaotic structures struggle to systematically navigate the vast design space of HDL, a limitation that causes them to converge prematurely on local optima. The failure of current paradigms to both effectively capitalize on valuable findings and adequately explore the solution space highlights the critical need for a more advanced framework.

To overcome these fundamental deficiencies, we propose **VERIMOA (Quality-guided Multi-path Mixture-of-Agents for HDL Generation)**, a framework that systematically addresses both noise susceptibility and constrained reasoning space in current multi-agent approaches. Building upon the Mixture-of-Agents (MoA) paradigm (Wang et al., 2024a), which aggregates outputs from the preceding layer, we introduce two key innovations. While MoA's layer-wise aggregation partially mitigates error propagation compared to linear pipelines, it still suffers from cascaded dependencies where hallucination errors accumulate across layers. We address noise susceptibility through a **quality-guided caching mechanism** that fundamentally breaks cascaded dependencies by evaluating and caching intermediate HDL outputs from all previous layers, enabling deeper agents to selectively leverage the highest-quality code from the entire process rather than being confined to the immediately preceding layer. For constrained reasoning space, we introduce **multi-path generation strategies** within each layer, incorporating intermediate representation paths via C++ and Python that decompose specification-to-HDL trans-

lation into two-stage processes (specification → high-level code → HDL), substantially expanding the solution space at each layer.

These two innovations synergistically enhance HDL generation performance. The quality-guided caching ensures monotonic knowledge accumulation across layers, preventing valuable insights from being lost. The multi-path generation leverages LLMs' rich parametric knowledge in C++ and Python while introducing structured diversity across generation paths (direct HDL, C++-guided, Python-guided) for comprehensive solution space exploration. Additionally, our framework can integrate simulation-based self-refinement to further enhance code quality, culminating in superior outputs.

Our contributions are summarized as follows:

- **Quality-Guided Mixture-of-Agents Architecture**: We propose a novel multi-agent framework that combines quality-guided caching with layered progress, ensuring monotonic knowledge accumulation and preventing information loss across layers.
- **Multi-Path Generation with Intermediate Representations**: We introduce parallel generation paths incorporating C++ and Python as intermediate representations, enabling two-stage HDL synthesis that leverages LLMs' strengths in high-level languages while promoting solution diversity.
- **Comprehensive Empirical Validation**: Through experiments on VerilogEval 2.0 and RTLLM 2.0 benchmarks, we demonstrate VERIMOA achieves superior performance with 15–30% improvements in Pass@1 across diverse LLM backbones, validating the effectiveness of our quality-guided multi-path approach.

## 2 PROBLEM DEFINITION

We formalize the HDL generation task and establish the optimization objective for benchmark evaluation. Given a hardware design benchmark $\mathcal{B}$ consisting of $N$ design problems, each problem $i \in \{1, 2, \ldots, N\}$ is characterized by a tuple $(\mathcal{D}_i, \mathcal{T}_i, \mathcal{S}_i)$, where $\mathcal{D}_i$ denotes the natural language design description, $\mathcal{T}_i$ represents the golden testbench for simulation-based verification, and $\mathcal{S}_i$ denotes the simulator that validates functional correctness. The objective is to construct an LLM-based generation system $\mathcal{G}$ that produces HDL code $H_i$ for each design description $\mathcal{D}_i$:

$$H_i = \mathcal{G}(\mathcal{D}_i; \theta) \tag{1}$$

where $\theta$ represents the system setting and generation strategy. The generated code $H_i$ is evaluated through simulation $\text{Pass}_i = \mathcal{S}_i(H_i, \mathcal{T}_i) \in \{\text{True}, \text{False}\}$, where $\text{Pass}_i = \text{True}$ indicates that $H_i$ passes all test cases, demonstrating both syntactic validity and functional correctness.
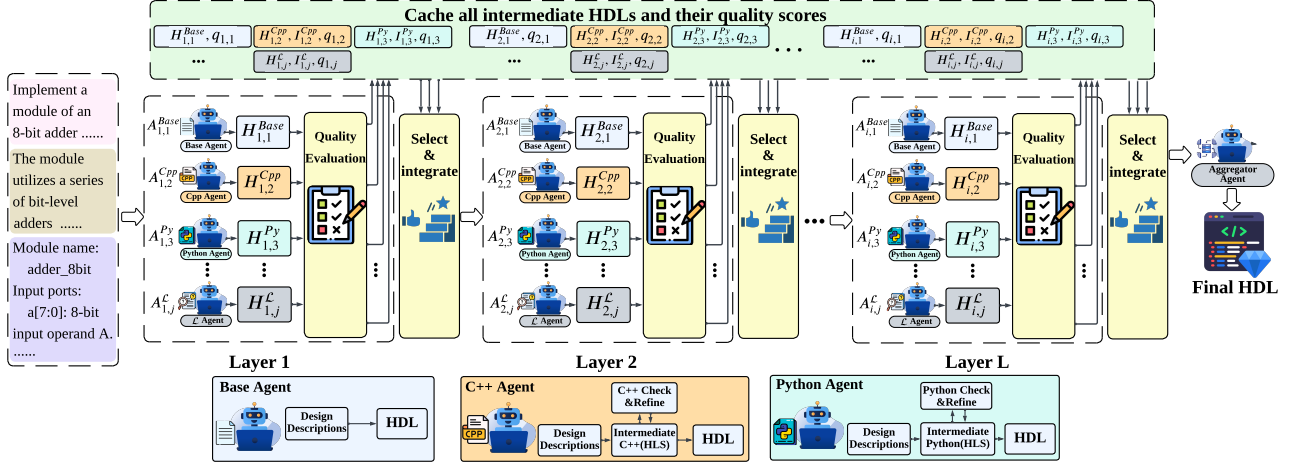
*Figure 1.* Overview of VERIMOA. The framework employs L layers with M agents per layer, featuring three agent types (Base, C++, Python) that generate HDL through different paths. A global cache stores all intermediate outputs with quality scores, enabling quality-guided selection across layers. Agents at each layer receive top-n quality HDL codes from all previous layers, ensuring monotonic improvement. The final aggregator synthesizes high-quality candidates into the output.

For practical deployment, the system generates multiple trials for each problem, and we measure performance using standard code generation metrics. Let $\mathcal{M}_k(\mathcal{G}, \mathcal{B})$ denote the pass@k metric that evaluates whether at least one correct solution exists among $k$ generated trials, averaged across all problems in benchmark $\mathcal{B}$. Our goal is to design a generation system $\mathcal{G}^*$ that maximizes this metric:

$$\mathcal{G}^* = \arg\max_{\mathcal{G}} \mathcal{M}_k(\mathcal{G}, \mathcal{B}) \qquad (2)$$

## 3  QUALITY-GUIDED MULTI-PATH MIXTURE-OF-AGENTS FOR HDL GENERATION

Existing multi-agent frameworks for HDL generation suffer from two critical deficiencies: susceptibility to noise propagation and constrained reasoning space exploration. To address these limitations, we propose **Quality-guided Multi-path Mixture-of-Agents for HDL Generation (VERIMOA)**, a framework that introduces two synergistic innovations: (1) a quality-guided caching mechanism ensuring monotonic knowledge accumulation across layers, and (2) multi-path generation strategies leveraging high-level code as intermediate representations for expanded solution space exploration. This section presents the implementation details of our framework, organized into two subsections corresponding to these core innovations.

### 3.1  Quality-Guided Architecture with Global Caching

As shown in Fig. 1, our quality-guided MoA consists of three interconnected components: MoA layers, a quality evaluator, and a global cache. The architecture is designed

to break the cascaded dependencies of standard MoA while preserving structured information flow across layers.

**MoA Layers.** The framework employs a hierarchical structure comprising $L$ layers, organized into proposer layers (layers 1 through $L-1$) and an aggregator layer (layer $L$). Each proposer layer contains $M$ agents operating in parallel to generate diverse HDL candidates. For an agent $A_{i,j}$ in the $i$-th proposer layer ($i \in \{1, \ldots, L-1\}$, $j \in \{1, \ldots, M\}$), the generation process is formalized as:

$$H_{i,j} = A_{i,j}(\mathcal{P}_{i,j}) \qquad (3)$$

where $\mathcal{P}_{i,j}$ denotes the prompt and $H_{i,j}$ represents the HDL output. Following standard MoA conventions (Wang et al., 2024a), we maintain consistent agent configurations across proposer layers ($A_{i,j} = A_{i',j}$ for all $i, i' \in \{1, \ldots, L-1\}$). The aggregator synthesizes high-quality candidates from previous layers into the final HDL output.

**Quality Evaluation.** The quality evaluator $\mathcal{Q}$ assesses each generated HDL candidate through simulation-based evaluation, producing a quality score $q_{i,j} = \mathcal{Q}(H_{i,j}, \mathcal{T}, \mathcal{S})$. This score reflects both syntactic validity and functional correctness, serving as the foundation for quality-guided selection. The evaluator implements a hierarchical scoring strategy detailed in Algorithm 1, which incorporates HDL domain-specific knowledge through three evaluation paths: (1) perfect scores for fully correct implementations, (2) severity-weighted penalties for functional failures with valid syntax, and (3) rule-based fallback evaluation for syntactic failures. The quality scores encode hardware design principles such as reset signal requirements, signal driving conflicts and timing constraints into explicit evaluation

**Algorithm 1** Quality Evaluation with HDL-Specific Criteria

**Require:** HDL code $H$, testbench $\mathcal{T}$, simulator $\mathcal{S}$

**Ensure:** Quality score $q \in [0, q_{\text{perfect}}]$

1: $\text{pass}_{\text{syntax}} \leftarrow \text{SyntaxTest}(H, \mathcal{S})$
2: $\text{pass}_{\text{func}} \leftarrow \text{FunctionTest}(H, \mathcal{T}, \mathcal{S})$
3: **if** $\text{pass}_{\text{syntax}} = \text{True}$ **and** $\text{pass}_{\text{func}} = \text{True}$ **then**
4:     **return** $q = q_{\text{perfect}}$
5: **else if** $\text{pass}_{\text{syntax}} = \text{True}$ **and** $\text{pass}_{\text{func}} = \text{False}$ **then**
6:     $p_{\text{severe}} \leftarrow \text{EvaluateSevereErrors}(H)$ {Logic errors, $\in [0, \alpha_{\text{severe}}]$}
7:     $p_{\text{moderate}} \leftarrow \text{EvaluateSynthesisIssues}(H)$ {Synthesis concerns, $\in [0, \alpha_{\text{moderate}}]$}
8:     $p_{\text{minor}} \leftarrow \text{EvaluateStyleIssues}(H)$ {Style issues, $\in [0, \alpha_{\text{minor}}]$}
9:     **return** $q = q_{\text{base}} - p_{\text{severe}} - p_{\text{moderate}} - p_{\text{minor}}$
10: **else**
11:     $q_{\text{structure}} \leftarrow \text{CheckModuleStructure}(H)$ {Module/IO, $\in [0, \beta_{\text{structure}}]$}
12:     $q_{\text{logic}} \leftarrow \text{CheckLogicKeywords}(H)$ {Logic constructs, $\in [0, \beta_{\text{logic}}]$}
13:     $q_{\text{format}} \leftarrow \text{CheckFormatting}(H)$ {Formatting, $\in [0, \beta_{\text{format}}]$}
14:     **return** $q = q_{\text{structure}} + q_{\text{logic}} + q_{\text{format}}$
15: **end if**

criteria, providing agents with HDL-specific feedback that injects domain expertise lacking in general-purpose LLMs.

**Global Cache and Quality-Guided Selection.** To address the unstable information propagation in standard MoA, we propose a global cache that maintains all intermediate outputs and enables quality-based selection across the entire generation process. All intermediate HDL outputs and their corresponding quality scores are stored in a global cache:

$$\mathcal{C} = \{(H_{i,j}, q_{i,j}) \mid i \in \{1, \ldots, L-1\}, j \in \{1, \ldots, M\}\} \quad (4)$$

This cache decouples information flow from rigid layer-by-layer dependencies, enabling agents to access high-quality references from any previous layer. For the first proposer layer, agents receive only the design description: $\mathcal{P}_{1,j} = \mathcal{D}$. For subsequent layers ($i \geq 2$), each agent receives an augmented prompt:

$$\mathcal{P}_{i,j} = \mathcal{D} \oplus \mathcal{H}_i^{(n)} \quad (5)$$

where $\oplus$ denotes concatenation and $\mathcal{H}_i^{(n)} = \text{TopN}(\mathcal{C}_{<i}, n)$ represents the top-$n$ highest-quality HDL codes selected from all previous layers. The TopN operation is defined as:

$$\text{TopN}(\mathcal{C}_{<i}, n) = \{H \mid (H, q) \in \mathcal{C}_{<i}, q \in \text{top-}n \text{ scores}\} \quad (6)$$

where $\mathcal{C}_{<i} = \{(H_{k,j}, q_{k,j}) \mid k < i, j \in \{1, \ldots, M\}\}$ denotes the cache subset containing outputs from all layers before layer $i$. This quality-guided selection ensures

that deeper agents always access the best available references, thereby breaking cascaded dependencies and enabling monotonic knowledge accumulation.

## 3.2 Multi-Path Generation with High-level Code as Intermediate Representations

While the quality-guided caching mechanism ensures monotonic knowledge accumulation at the architectural level, we further enhance the framework by improving proposer layers internally through the multi-path generation strategy. The strategy leverages high-level programming languages as intermediate representations, enabling LLMs to utilize their superior C++/Python knowledge while exploring diverse reasoning paths within each layer. Let $\mathcal{L} \in \{\text{C++}, \text{Python}\}$ denote the intermediate language. For an agent following the $\mathcal{L}$-path in layer $i$, generation proceeds in two stages:

**Stage 1: Specification to Intermediate Code.** To further leverage the quality-guided caching mechanism, we evaluate and cache intermediate codes analogously to HDL codes, storing them in intermediate code cache $\mathcal{C}_{\mathcal{L}} = \{(I_{i,j}^{\mathcal{L}}, q_{i,j}^{\mathcal{L}}) \mid i \in \{1, \ldots, L-1\}, j \in \{1, \ldots, M_{\mathcal{L}}\}\}$ for each language $\mathcal{L}$. The agent generates intermediate code as:

$$I_{i,j}^{\mathcal{L}} = A_{i,j}^{\text{Stage1}}(\mathcal{P}_{i,j}^{\mathcal{L}}) \quad (7)$$

where the prompt is constructed as:

$$\mathcal{P}_{i,j}^{\mathcal{L}} = \begin{cases} \mathcal{D} & \text{if } i = 1 \\ \mathcal{D} \oplus \mathcal{I}_i^{\mathcal{L},(k)} & \text{if } i \geq 2 \end{cases} \quad (8)$$

For the first layer, the agent generates intermediate code directly from the design description. For subsequent layers ($i \geq 2$), the prompt is augmented with $\mathcal{I}_i^{\mathcal{L},(k)} = \text{TopK}(\mathcal{C}_{\mathcal{L}, <i}, k)$, which represents the top-$k$ highest-quality intermediate codes selected from all previous layers. We then perform syntax validation using checking scripts and invoke the agent for self-refinement based on cheching feedback, producing refined intermediate code $I_{i,j}^{\mathcal{L},\text{refined}}$.

**Stage 2: Intermediate Code to HDL.** Using the refined intermediate code as a reference, the agent generates HDL:

$$H_{i,j}^{\mathcal{L}} = A_{i,j}^{\text{Stage2}}(\mathcal{D} \oplus I_{i,j}^{\mathcal{L},\text{refined}} \oplus \mathcal{H}_i^{(n)}) \quad (9)$$

where $\mathcal{H}_i^{(n)}$ denotes the top-$n$ HDL references from the global cache. For $i = 1$, this simplifies to $H_{1,j}^{\mathcal{L}} = A_{1,j}^{\text{Stage2}}(\mathcal{D} \oplus I_{1,j}^{\mathcal{L},\text{refined}})$ as no prior HDL references exist.

To enable quality-guided selection of intermediate codes across layers, we assign each intermediate code a quality score based on the HDL it generates:

$$q_{i,j}^{\mathcal{L}} = q(H_{i,j}^{\mathcal{L}}) \quad (10)$$

where $q(H_{i,j}^{\mathcal{L}})$ is evaluated by the evaluator $\mathcal{Q}$ described in Algorithm 1. This task-aligned evaluation prioritizes intermediate codes that enable high-quality HDL translation.

**Multi-Path Proposer Configuration.** We define three agents: **Baseline agents** ($A_{i,j}^{\text{Base}}$) generate HDL directly from the design description, **C++ agents** ($A_{i,j}^{\text{Cpp}}$) and **Python agents** ($A_{i,j}^{\text{Py}}$) follow the two-stage process with C++ and Python as intermediate languages, respectively. Each proposer layer consists of $M$ agents operating in parallel, with each position $j \in \{1, \ldots, M\}$ assigned one agent type:

$$\mathcal{A}_i = \{A_{i,j}^{t_j} \mid j \in \{1, \ldots, M\}, t_j \in \{\text{Base}, \text{Cpp}, \text{Py}\}\} \tag{11}$$

This heterogeneous mixture explores multiple reasoning paths simultaneously, where baseline agents prioritize hardware idioms, C++ agents emphasize bit-level control, while Python agents leverage high-level expressiveness, leading to diverse HDL implementations for identical specifications.

**Optional Simulator-Based Self-Refinement.** Following established practices (Mi et al., 2024; Zhao et al., 2024), our framework can optionally integrate simulator-based self-refinement. After generating initial HDL output $H_{i,j}$, agents invoke simulator $\mathcal{S}$ with testbench $\mathcal{T}$ to obtain execution feedback for refinement. This mechanism applies uniformly to all agent types in both proposer and aggregator layers.

# 4 THEORETICAL ANALYSIS OF QUALITY-GUIDED MULTI-PATH GENERATION

This section provides theoretical analysis of our VERIMOA framework, establishing formal guarantees for performance improvement and explaining how quality-guided caching and multi-path generation enhance HDL generation. We connect our design to established Mixture-of-Agents principles, prove monotonic quality guarantees, and analyze how multi-path strategies amplify both quality and diversity.

## 4.1 Foundation: MoA Performance Decomposition

Recent empirical analysis of Mixture-of-Agents (Li et al., 2025) demonstrates that MoA performance exhibits strong positive correlation with two key factors: proposer quality and proposer diversity. Through statistical analysis across reasoning benchmarks, they establish that MoA performance $t$ can be approximated by a linear model:

$$t = \alpha \times q + \beta \times d + \gamma \tag{12}$$

where $q$ represents proposer quality, $d$ represents proposer diversity, and $\alpha, \beta, \gamma$ are regression coefficients with $\alpha > \beta$, indicating quality exerts stronger influence than diversity.

In HDL generation, proposer quality encompasses agent capability and reference code quality from previous layers. Let $\mathcal{A}_i = \{A_{i,1}, \ldots, A_{i,M}\}$ denote the agent configuration at layer $i$, and $\mathcal{H}_i^{(n)}$ the reference HDL set. The effective proposer quality is:

$$Q_i = f(\text{capability}(\mathcal{A}_i), \text{quality}(\mathcal{H}_i^{(n)})) \tag{13}$$

where $f$ captures the combined effect of agent capability and reference quality. With fixed agent configurations across layers, improving $Q_i$ reduces to enhancing quality($\mathcal{H}_i^{(n)}$). Given $\alpha > \beta$, our framework prioritizes maximizing reference quality through quality-guided caching while promoting diversity through multi-path generation.

## 4.2 Quality Guarantee through Global Caching

**Error Accumulation in Standard MoA.** Applying standard MoA to HDL generation causes performance degradation as layer depth increases. This occurs because LLMs' limited HDL knowledge causes each layer to introduce non-trivial errors. In standard MoA, each proposer layer only receives outputs from the preceding layer without quality-based selection. Consequently, low-quality outputs from layer $i - 1$ contaminate layer $i$'s reference set, harming final performance through cascaded error propagation.

**Monotonic Quality Improvement.** Our quality-guided caching mechanism addresses error accumulation by maintaining a global cache and selecting only top-$n$ quality outputs across all previous layers, ensuring monotonic improvement while breaking rigid layer-by-layer propagation.

For any layer $i \geq 2$, the quality guarantee can be formalized as:

$$\min_{H \in \mathcal{H}_{i+1}^{(n)}} q(H) \geq \min_{H \in \mathcal{H}_i^{(n)}} q(H) \tag{14}$$

$$\frac{1}{n} \sum_{H \in \mathcal{H}_{i+1}^{(n)}} q(H) \geq \frac{1}{n} \sum_{H \in \mathcal{H}_i^{(n)}} q(H) \tag{15}$$

Equation (14) ensures the minimum quality is non-decreasing, while Equation (15) guarantees average quality improvement. These properties hold by construction of the TopN operation (Equation (6)): since $\mathcal{C}_{<i+1} \supset \mathcal{C}_{<i}$, selecting top-$n$ from a strictly larger set that includes all previous top-$n$ elements plus new candidates from layer $i$ cannot decrease either minimum or average quality.

This quality guarantee establishes non-decreasing information quality flow across proposer layers. Consequently, agents in deeper layers operate on strictly superior reference materials, leading to expected quality improvements:

$$\mathbb{E}[q_{i+1,j}] \geq \mathbb{E}[q_{i,j}], \quad \forall j \in \{1, \ldots, M\} \tag{16}$$

where the expectation is taken over LLM generation stochasticity. This monotonicity property prevents information loss during multi-layer process, enabling consistent performance gains as depth increases, in contrast to the degradation observed in standard MoA.

## 4.3 Synergistic Enhancement through Multi-Path Generation

The multi-path generation strategy synergistically addresses both dimensions of MoA performance (Equation (12)): proposer quality and diversity.

**Quality Enhancement via Intermediate Representations.** Recall that proposer quality $Q_i$ depends on both agent capability and reference code quality (Equation (13)). The two-stage mechanism enhances agent capability by leveraging LLMs' substantially stronger parametric knowledge of C++/Python (Xiong et al., 2024). In Stage 1, the LLM decomposes abstract design requirements into explicit algorithmic constructs (control flow, data structures, and computational patterns) that directly correspond to hardware behavior (Xu et al., 2024). This intermediate representation serves as a structured scaffold for Stage 2 HDL generation, effectively translating the LLM's strong high-level language comprehension into accurate hardware implementations.

To formalize quality propagation in multi-path generation, consider an agent following the $\mathcal{L}$-path at layer $i$. The expected HDL quality depends on both intermediate code quality and reference HDL quality:

$$\mathbb{E}[q(H_{i,j}^{\mathcal{L}})] = \mathbb{E}[g(q(I_{i,j}^{\mathcal{L}}), \text{quality}(\mathcal{H}_i^{(n)}))] \quad (17)$$

where $g$ represents the Stage 2 translation process, whose expected output quality increases monotonically with both input arguments. Quality-guided caching for intermediate codes ensures:

$$\min_{I \in \mathcal{I}_{i+1}^{\mathcal{L},(k)}} q(I) \geq \min_{I \in \mathcal{I}_i^{\mathcal{L},(k)}} q(I) \quad (18)$$

$$\frac{1}{k} \sum_{I \in \mathcal{I}_{i+1}^{\mathcal{L},(k)}} q(I) \geq \frac{1}{k} \sum_{I \in \mathcal{I}_i^{\mathcal{L},(k)}} q(I) \quad (19)$$

Combined with HDL quality guarantees (Equations (14) (15)), we have:

$$\text{quality}(\mathcal{H}_{i+1}^{(n)}) \geq \text{quality}(\mathcal{H}_i^{(n)}) \quad (20)$$

Since $g$ is expected to be monotonically increasing and both inputs improve, applying the monotonicity yields:

$$\begin{aligned} \mathbb{E}[q(H_{i+1,j}^{\mathcal{L}})] &= \mathbb{E}[g(q(I_{i+1,j}^{\mathcal{L}}), \text{quality}(\mathcal{H}_{i+1}^{(n)}))] \\ &\geq \mathbb{E}[g(q(I_{i,j}^{\mathcal{L}}), \text{quality}(\mathcal{H}_i^{(n)}))] = \mathbb{E}[q(H_{i,j}^{\mathcal{L}})] \end{aligned} \quad (21)$$

This monotonicity extends across all agent types (Baseline, C++, Python), collectively ensuring $Q_{i+1} \geq Q_i$.

**Diversity Amplification through Heterogeneous Paths.** The multi-path architecture amplifies diversity through fundamentally different reasoning trajectories: baseline agents translate directly to HDL, C++ agents reason through imperative abstractions, while Python agents leverage high-level expressiveness. Let $\mathcal{H}_i^{\text{Base}}$, $\mathcal{H}_i^{\text{Cpp}}$, and $\mathcal{H}_i^{\text{Py}}$ denote the HDL outputs from each agent type at layer $i$. The structural diversity is characterized by:

$$d_i = \mathbb{E}[\text{dissimilarity}(\mathcal{H}_i^{\text{Base}}, \mathcal{H}_i^{\text{Cpp}}, \mathcal{H}_i^{\text{Py}})] \quad (22)$$

where dissimilarity captures differences in code structure and implementation patterns, directly enhancing the diversity term $d$ in Equation (12).

**Synergistic Performance Maximization.** The heterogeneous mixture of quality-improving paths (Equation (21)), each exploring distinct solution spaces (Equation (22)), synergistically maximizes MoA performance according to Equation (12). Quality-guided caching ensures monotonic accumulation of high-quality references across all paths, while multi-path generation expands the solution space through diverse reasoning strategies.

## 5 EXPERIMENTS

We conduct comprehensive experiments to evaluate the effectiveness of our VERIMOA framework. Our evaluation is designed to address four key research questions:

**RQ1:** How does VERIMOA compare to state-of-the-art methods in HDL generation performance?

**RQ2:** What is the contribution of each core component to the overall performance?

**RQ3:** How sensitive is the framework to variations in key parameters?

**RQ4:** Can VERIMOA effectively explore diverse solutions and propagate high-quality information across layers?

### 5.1 Experimental Setup

**Benchmarks.** We evaluate on VerilogEval 2.0 (Pinckney et al., 2025) (156 problems covering combinational and sequential circuits) and RTLLM 2.0 (Liu et al., 2024b) (50 complex real-world design tasks).

**Baselines.** We compare VERIMOA against representative methods from two categories. For non-training approaches, we include Direct (direct generation without prompting strategies), CoT (Wei et al., 2022) (Chain-of-Thought prompting), HDLCoRe (Ping et al., 2025) (difficulty-adaptive prompting with RAG), and VeriMaAS (Bhattaram et al., 2025) (adaptive multi-agent system), evaluated across various LLMs including Qwen2.5 series and GPT-4o models. For fine-tuned models, we include RTLCoder-Mistral (Liu et al., 2024a), RTLCoder-DeepSeek-Coder (Liu et al., 2024a), OriGen-DeepSeek-Coder-7B-v1.5 (Cui et al., 2024), HaVen-CodeQwen1.5 (Yang et al.,

*Table 1.* Performance comparison of non-training baselines and our VERIMOA on VerilogEval 2.0 and RTLLM 2.0 benchmarks.

| Model | Method | VerilogEval 2.0 (%) | | | RTLLM 2.0 (%) | | |
|---|---|---|---|---|---|---|---|
| | | *Pass@1* | *Pass@3* | *Pass@5* | *Pass@1* | *Pass@3* | *Pass@5* |
| Qwen2.5-7B | Direct | 22.90 | 32.20 | 40.39 | 18.99 | 30.53 | 36.51 |
| | CoT | 26.87 ↑3.97 | 35.33 ↑3.13 | 42.58 ↑2.19 | 27.49 ↑8.50 | 34.76 ↑4.23 | 39.95 ↑3.44 |
| | HDLCoRe | 30.92 ↑8.02 | 37.86 ↑5.66 | 46.49 ↑6.10 | 34.68 ↑15.69 | 40.32 ↑9.79 | 44.82 ↑8.31 |
| | VeriMaAS | 32.81 ↑9.91 | 40.25 ↑8.05 | 49.07 ↑8.68 | 40.33 ↑21.34 | 47.24 ↑16.71 | 50.47 ↑13.96 |
| | **VERIMOA** | **56.44 ↑33.54** | **62.07 ↑29.87** | **65.27 ↑24.88** | **52.07 ↑33.08** | **58.10 ↑27.57** | **60.46 ↑23.95** |
| Qwen2.5-Coder-7B | Direct | 33.91 | 43.90 | 48.06 | 26.16 | 36.94 | 41.11 |
| | CoT | 36.74 ↑2.83 | 45.88 ↑1.98 | 50.99 ↑2.93 | 31.87 ↑5.71 | 39.36 ↑2.42 | 45.36 ↑4.25 |
| | HDLCoRe | 40.67 ↑6.76 | 47.81 ↑3.91 | 52.64 ↑4.58 | 36.26 ↑10.10 | 45.61 ↑8.67 | 51.12 ↑10.01 |
| | VeriMaAS | 44.73 ↑10.82 | 51.62 ↑7.72 | 55.78 ↑7.72 | 44.71 ↑18.55 | 51.37 ↑14.43 | 56.46 ↑15.35 |
| | **VERIMOA** | **60.96 ↑27.05** | **68.13 ↑24.23** | **70.69 ↑22.63** | **54.43 ↑28.27** | **62.52 ↑25.58** | **66.24 ↑25.13** |
| Qwen2.5-14B | Direct | 39.26 | 49.98 | 53.74 | 31.71 | 40.20 | 45.87 |
| | CoT | 41.24 ↑1.98 | 51.62 ↑1.64 | 54.37 ↑0.63 | 36.49 ↑4.78 | 44.64 ↑4.44 | 48.29 ↑2.42 |
| | HDLCoRe | 45.18 ↑5.92 | 55.83 ↑5.85 | 59.42 ↑5.68 | 41.78 ↑10.07 | 49.97 ↑9.77 | 53.17 ↑7.30 |
| | VeriMaAS | 48.97 ↑9.71 | 57.84 ↑7.86 | 61.77 ↑8.03 | 48.14 ↑16.43 | 54.28 ↑14.08 | 58.76 ↑12.89 |
| | **VERIMOA** | **65.27 ↑26.01** | **71.61 ↑21.63** | **74.64 ↑20.90** | **55.06 ↑23.35** | **62.76 ↑22.56** | **66.33 ↑20.46** |
| Qwen2.5-Coder-14B | Direct | 39.74 | 49.72 | 52.21 | 35.61 | 42.34 | 46.43 |
| | CoT | 43.81 ↑4.07 | 51.27 ↑1.55 | 54.11 ↑1.90 | 40.65 ↑5.04 | 48.77 ↑6.43 | 51.72 ↑5.29 |
| | HDLCoRe | 46.82 ↑7.08 | 54.79 ↑5.07 | 57.04 ↑4.83 | 47.86 ↑12.25 | 55.19 ↑12.85 | 58.18 ↑11.75 |
| | VeriMaAS | 50.96 ↑11.22 | 58.48 ↑8.76 | 62.63 ↑10.42 | 51.20 ↑15.59 | 58.93 ↑16.59 | 61.67 ↑15.24 |
| | **VERIMOA** | **66.86 ↑27.12** | **73.24 ↑23.52** | **76.38 ↑24.17** | **59.76 ↑24.15** | **64.33 ↑21.99** | **67.22 ↑20.79** |
| Qwen2.5-32B | Direct | 46.85 | 56.11 | 58.41 | 43.62 | 48.73 | 50.75 |
| | CoT | 48.72 ↑1.87 | 57.25 ↑1.14 | 59.80 ↑1.39 | 46.99 ↑3.37 | 50.50 ↑1.77 | 52.46 ↑1.71 |
| | HDLCoRe | 51.75 ↑4.90 | 59.63 ↑3.52 | 61.79 ↑3.38 | 50.95 ↑7.33 | 54.06 ↑5.33 | 58.98 ↑8.23 |
| | VeriMaAS | 53.57 ↑6.72 | 61.82 ↑5.71 | 64.26 ↑5.85 | 54.18 ↑10.56 | 57.21 ↑8.48 | 61.54 ↑10.79 |
| | **VERIMOA** | **71.85 ↑25.00** | **76.48 ↑20.37** | **78.16 ↑19.75** | **63.81 ↑20.19** | **67.60 ↑18.87** | **69.88 ↑19.13** |
| Qwen2.5-Coder-32B | Direct | 46.93 | 55.73 | 57.12 | 40.40 | 45.42 | 48.24 |
| | CoT | 48.65 ↑1.72 | 56.31 ↑0.58 | 59.34 ↑2.22 | 45.72 ↑5.32 | 48.75 ↑3.33 | 50.29 ↑2.05 |
| | HDLCoRe | 51.28 ↑4.35 | 58.63 ↑2.90 | 61.47 ↑4.35 | 51.72 ↑11.32 | 54.73 ↑9.31 | 59.64 ↑11.40 |
| | VeriMaAS | 56.67 ↑9.74 | 63.46 ↑7.73 | 66.92 ↑9.80 | 55.82 ↑15.42 | 59.97 ↑14.55 | 62.31 ↑14.07 |
| | **VERIMOA** | **73.31 ↑26.38** | **79.05 ↑23.32** | **81.27 ↑24.15** | **65.49 ↑25.09** | **71.42 ↑26.00** | **74.11 ↑25.87** |
| GPT-4o-mini | Direct | 48.97 | 56.94 | 58.62 | 46.60 | 50.71 | 51.95 |
| | CoT | 52.20 ↑3.23 | 59.83 ↑2.89 | 60.93 ↑2.31 | 48.27 ↑1.67 | 53.56 ↑2.85 | 55.48 ↑3.53 |
| | HDLCoRe | 53.50 ↑4.53 | 61.46 ↑4.52 | 63.72 ↑5.10 | 51.19 ↑4.59 | 56.72 ↑6.01 | 59.01 ↑7.06 |
| | VeriMaAS | 57.24 ↑8.27 | 64.77 ↑7.83 | 66.85 ↑8.23 | 57.25 ↑10.65 | 61.46 ↑10.75 | 63.17 ↑11.22 |
| | **VERIMOA** | **72.43 ↑23.46** | **76.94 ↑20.00** | **79.46 ↑20.84** | **64.23 ↑17.63** | **67.45 ↑16.74** | **68.67 ↑16.72** |
| GPT-4o | Direct | 64.74 | 69.89 | 71.66 | 52.48 | 56.18 | 57.62 |
| | CoT | 66.38 ↑1.64 | 71.18 ↑1.29 | 74.22 ↑2.56 | 54.89 ↑2.41 | 59.22 ↑3.04 | 61.77 ↑4.15 |
| | HDLCoRe | 69.60 ↑4.86 | 73.53 ↑3.64 | 75.86 ↑4.20 | 56.27 ↑3.79 | 61.47 ↑5.29 | 63.23 ↑5.61 |
| | VeriMaAS | 71.34 ↑6.60 | 76.12 ↑6.23 | 79.21 ↑7.55 | 61.70 ↑9.22 | 67.25 ↑11.07 | 68.84 ↑11.22 |
| | **VERIMOA** | **84.97 ↑20.23** | **89.65 ↑19.76** | **91.03 ↑19.37** | **69.17 ↑16.69** | **73.68 ↑17.50** | **75.62 ↑18.00** |

2025), VeriRL-DeepSeek-Coder (Teng et al., 2025), and VeriRL-CodeQwen2.5 (Teng et al., 2025).

**Evaluation Metric.** We evaluate HDL generation performance from the perspective of functional correctness using the widely adopted pass@k metric (Pinckney et al., 2025), which estimates the proportion of problems for which at least one of $k$ generated solutions passes functional verification. The metric is defined as:

$$\text{pass@}k := \mathbb{E}\left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}\right] \quad (23)$$

where $n \geq k$ is the number of sampled outputs per problem and $c$ is the number of correct outputs among them. Following the experimental setting in VerilogEval 2.0 (Pinckney et al., 2025), we set $n = 10$ in all our experiments.

**Implementation Details.** We use Icarus Verilog (iverilog)

as our simulator for its efficiency and ease of integration. For LLM generation, we adopt the sampling strategy with temperature=0.8 and top_p=0.95, following the configuration in VerilogEval 2.0 (Pinckney et al., 2025) to ensure fair comparison. Our VERIMOA framework uses 4 proposer layers with 6 parallel agents per layer, consistent with the standard MoA configuration (Wang et al., 2024a). The code is available at: https://github.com/hping666/HDL_Generation.

## 5.2 HDL Generation Performance (RQ1)

We evaluate VERIMOA's HDL generation performance against state-of-the-art baselines on VerilogEval 2.0 and RTLLM 2.0 benchmarks. Tables 1 and 2 present comprehensive results across multiple LLM backbones and baseline methods.

*Table 2.* Performance comparison of fine-tuned Verilog-specific models and our VERIMOA on VerilogEval 2.0 and RTLLM 2.0 benchmarks. First , second , and third best results are highlighted for each metric.

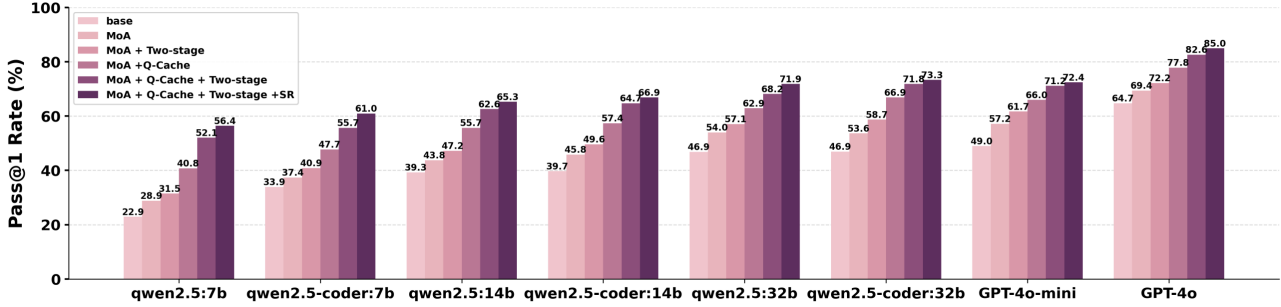| Model | | Size | VerilogEval 2.0 (%) | | | RTLLM 2.0 (%) | | |
|---|---|---|---|---|---|---|---|---|
| | | | *Pass@1* | *Pass@3* | *Pass@5* | *Pass@1* | *Pass@3* | *Pass@5* |
| Verilog-Specific (finetune) | RTLCoder-Mistral | 7B | 35.62 | 36.63 | 37.71 | 38.68 | 40.86 | 41.84 |
| | RTLCoder-DeepSeek-Coder | 6.7B | 39.67 | 43.66 | 45.99 | 40.75 | 46.94 | 49.83 |
| | OriGen-DeepSeek-Coder-7B-v1.5 | 7B | 51.19 | 54.88 | 56.78 | 41.11 | 52.09 | 58.48 |
| | HaVen-CodeQwen1.5 | 7B | 55.40 | 59.73 | 62.76 | 51.04 | 57.38 | 60.89 |
| | VeriRL-DeepSeek-Coder | 6.7B | 64.57 | 68.96 | 71.78 | 58.64 | 63.92 | 66.26 |
| | VeriRL-CodeQwen2.5 | 7B | 66.28 | 71.35 | 73.43 | 61.53 | 65.27 | 68.94 |
| Ours | Qwen2.5-Coder-7B + VERIMOA | 7B | 60.96 | 68.13 | 70.69 | 54.43 | 62.52 | 66.24 |
| | Qwen2.5-Coder-14B + VERIMOA | 14B | 66.86 | 73.24 | 76.38 | 59.76 | 64.33 | 66.22 |
| | Qwen2.5-Coder-32B + VERIMOA | 32B | 73.31 | 79.05 | 81.27 | 65.49 | 71.42 | 74.11 |



*Figure 2.* Ablation Study of techniques adopted in our framework. Pass@1 performance [Higher is Better] reported on VerilogEval 2.0 dataset.

**Comparison with Non-Training Methods.** Table 1 demonstrates that VERIMOA consistently outperforms all non-training baselines across different LLM backbones and both benchmarks, with particularly significant improvements on smaller LLMs. On VerilogEval 2.0, VERIMOA with Qwen2.5-7B achieves 56.44% Pass@1, outperforming Direct (22.90%) by 33.54 points and VeriMaAS (32.81%) by 23.63 points. Remarkably, VERIMOA with Qwen2.5-7B (56.44%) surpasses VeriMaAS with the larger Qwen2.5-32B (53.57%), and VERIMOA with Qwen2.5-Coder-7B (60.96%) exceeds VeriMaAS with Qwen2.5-Coder-32B (56.67%), showing that our framework enables smaller LLMs to match and even exceed larger models through superior information propagation and solution space exploration. At larger scales with Qwen2.5-Coder-32B, VERIMOA achieves 73.31% versus VeriMaAS's 56.67% (16.64 points improvement). Performance gains extend to commercial models: GPT-4o-mini reaches 72.43% versus VeriMaAS's 56.24% (16.19 points improvement), and GPT-4o attains 84.97% versus 71.34% (13.63 points improvement). On RTLLM 2.0, VERIMOA consistently outperforms VeriMaAS by 7–12 points of Pass@1 across model scales, with stronger gains for smaller models. These results validate that our quality-guided caching and multi-path generation strategies enchance HDL generation performance through architectural innovations applicable to diverse LLMs.

**Comparison with Fine-Tuned Models.** Table 2 shows that VERIMOA achieves competitive or superior performance without requiring training. With Qwen2.5-Coder-7B, VERIMOA attains 60.96% Pass@1 on VerilogEval 2.0, approaching fine-tuned models like VeriRL-DeepSeek-Coder (64.57%) and demonstrating that effective architectural design can partially compensate for domain-specific training. With Qwen2.5-Coder-14B, VERIMOA achieves 66.86%, matching the best fine-tuned 7B model VeriRL-CodeQwen2.5 (66.28%). Using Qwen2.5-Coder-32B, VERIMOA reaches 73.31%, surpassing all fine-tuned models by substantial margins (7.03 points improvement over VeriRL-CodeQwen2.5). This demonstrates our training-free framework can match or exceed fine-tuned models through quality-guided progress and diverse solution exploration.

### 5.3 Ablation Study (RQ2)

We conduct ablation studies to analyze the contribution of each core component. We evaluate six configurations across multiple LLM backbones using Pass@1 on VerilogEval 2.0 basesd on the following techniques: **Base** (direct generation), **MoA** (standard Mixture-of-Agents), **Two-stage** (high-level code as intermediate representations), **Q-Cache** (quality-guided caching mechanism), and **SR** (simulator-based self-refinement). Figure 2 reveals that Q-Cache provides substantially greater benefits than Two-stage alone.
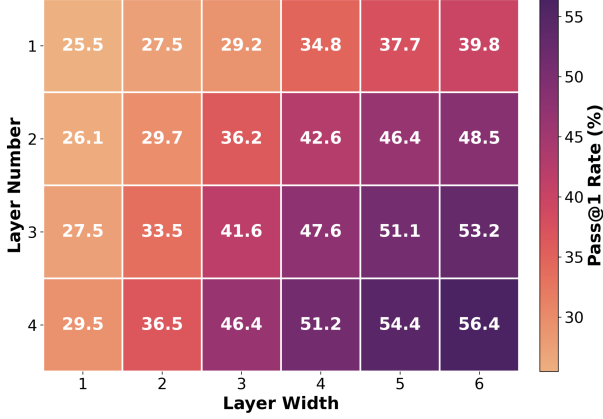
*Figure 3.* VeriMoA performance (Pass@1)[Darker is Better] on VerilogEval 2.0 using Qwen2.5:7b with varying layer configurations.

On Qwen2.5-7B, adding Q-Cache to MoA improves Pass@1 from 28.86% to 40.79% (11.93 points improvement), while adding Two-stage to MoA only yields 31.46% (2.60 points improvement). This pattern holds across all models: with GPT-4o, Q-Cache contributes 8.43 points improvement versus Two-stage's 2.78 points improvement.

More importantly, Q-Cache is essential for realizing Two-stage's full potential. MoA + Two-stage achieves modest gains, but MoA + Q-Cache + Two-stage shows substantial improvements. For Qwen2.5-7B, adding Two-stage on top of Q-Cache improves from 40.79% to 52.06% (11.27 points), while adding Two-stage to baseline MoA yields only 2.60 points improvement, demonstrating that Q-Cache enables effective multi-path exploration with high-quality information propagation. The complete framework with SR achieves the best performance across all models, with particularly strong gains for smaller LLMs (Qwen2.5-7B: 22.90% to 56.44%). These results validate that Q-Cache serves as the critical enabler for effective multi-path generation.

### 5.4 Parameter Sensitivity Analysis (RQ3)

We investigate the impact of two key hyperparameters on QMMoA's performance: layer number (depth) and layer width (number of parallel agents per layer). Figure 3 presents Pass@1 results on VerilogEval 2.0 using Qwen2.5-7B across different configurations.

The results reveal that both dimensions must reach sufficient scale for strong performance. With minimal width (1 agent) or shallow depth (1 layer), performance remains low (25.5%-39.8%). Significant improvements emerge only when both parameters reach adequate thresholds: configurations with 3+ layers and 4+ agents achieve 47.6%-56.4% Pass@1, substantially outperforming sparse configurations. This validates that our framework requires sufficient capac-
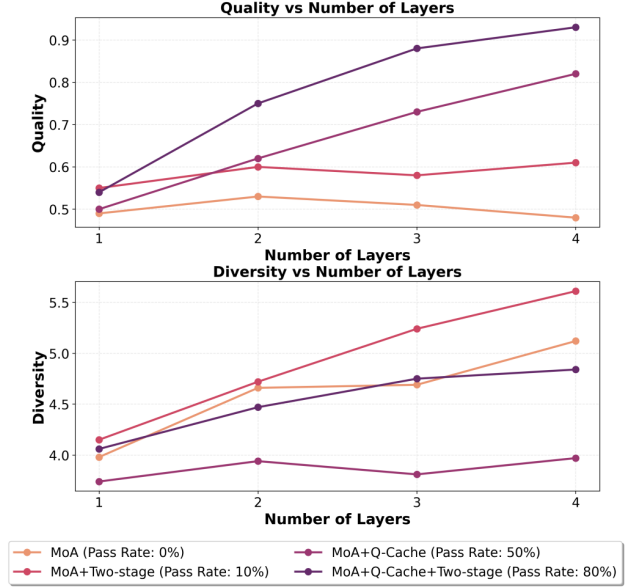


*Figure 4.* Quality and diversity metrics during LIFObuffer generation from RTLLM 2.0 across 10 trials.

ity in both dimensions to effectively accumulate high-quality solutions and explore diverse reasoning paths.

When comparing configurations with equal total agents, increasing width yields greater gains than increasing depth. For instance, with 12 total agents, the 3-layer-4-width configuration achieves 47.6% versus the 4-layer-3-width configuration's 46.4%. Furthermly, 2-layer-6-width achieves 48.5%, outperforming the former both. This demonstrates that width (parallel diversity) maters importantly, validating our multi-path generation strategy's emphasis on exploring diverse solution spaces through heterogeneous agent mixtures within each layer.

### 5.5 Case Study: Quality and Diversity Analysis (RQ4)

We analyze quality and diversity evolution on a complex task from RTLLM 2.0, named "LIFObuffer", by using Qwen2.5-7B with layer width of 6. We measure the top-6 cached HDL candidates at each layer across four configurations, averaging metrics over 10 trials. Quality is evaluated using our scoring mechanism, while diversity is measured by Vendi score (Li et al., 2025). Figure 4 reveals distinct patterns. Standard MoA and MoA+Two-stage show stagnant quality (0.49 to 0.48 and 0.55 to 0.61) with 0-10% pass rates. MoA+Q-Cache achieves monotonic quality improvement (0.50 to 0.82, 50% pass rate), while MoA+Q-Cache+Two-stage demonstrates the strongest growth (0.54 to 0.93, 80% pass rate). This validates that Q-Cache enables progressive quality improvement, which Two-stage amplifies when built upon this stable foundation.

For diversity, MoA and MoA+Two-stage maintain high values (3.98-5.61) but achieve low pass rates, indicating unfocused exploration. MoA+Q-Cache shows modest diversity (3.74-3.97) but reasonable performance through quality filtering. MoA+Q-Cache+Two-stage achieves both high diversity (4.06-4.84) and the best pass rate (80%), confirming that diversity alone is insufficient—effective exploration requires quality-guided selection to ensure diverse candidates contribute meaningful improvements rather than noise.

# 6 CONCLUSION

This paper introduces VERIMOA, a quality-guided multi-path Mixture-of-Agents framework that systematically addresses noise susceptibility and constrained reasoning space in multi-agent HDL generation. By combining a quality-guided caching mechanism that ensures monotonic knowledge accumulation with a multi-path generation strategy leveraging C++ and Python as intermediate representations, VERIMOA achieves 15–30% improvements in Pass@1 across diverse LLM backbones on both VerilogEval 2.0 and RTLLM 2.0 benchmarks. Our framework enables smaller LLMs to match and even exceed larger models, fine-tuned alternatives without costly training, demonstrating that effective architectural design through quality-guided progress and diverse solution exploration offers a practical and scalable approach for automated RTL design.

## REFERENCES

Alsaqer, S., Alajmi, S., Ahmad, I., and Alfailakawi, M. The potential of llms in hardware design. *Journal of Engineering Research*, 13(3):2392–2404, 2025.

Bhattaram, A., Ramamoorthy, J., Gupta, R., Marculescu, D., and Stamoulis, D. Automated multi-agent workflows for RTL design. In *39th Conference on Neural Information Processing Systems (NeurIPS 2025) Workshop: ML for Systems*, 2025.

Chen, Z., Chang, Y., Du, Z., Zhang, X., Kang, Z., Yu, Y., Bai, C., Huang, Y., and Chen, L. ChipSeek-R1: Generating human-surpassing RTL with LLM via hierarchical reward-driven reinforcement learning. *arXiv preprint arXiv:2507.04736*, 2025.

Cui, F., Yin, C., Zhou, K., Xiao, Y., Sun, G., Xu, Q., Guo, Q., Liang, Y., Zhang, X., Song, D., and Lin, D. Origen: Enhancing rtl code generation with code-to-code augmentation and self-reflection. In *Proceedings of the 43rd IEEE/ACM International Conference on Computer-Aided Design*, pp. 1–9. IEEE/ACM, 2024.

DeLorenzo, M., Chowdhury, A. B., Gohil, V., Thakur, S., Pearce, H., Karri, R., and Rajendran, J. Abstractions-of-thought: Intermediate representations for LLM reasoning in hardware design. *arXiv preprint arXiv:2505.15873*, 2025.

Gao, M., Zhao, J., Lin, Z., Ding, W., Hou, X., Feng, Y., Li, C., and Guo, M. AutoVCoder: A systematic framework for automated Verilog code generation using LLMs. *arXiv preprint arXiv:2407.18333*, 2024.

Li, W., Lin, Y., Xia, M., and Jin, C. Rethinking mixture-of-agents: Is mixing different large language models beneficial? *arXiv preprint arXiv:2502.00674*, 2025.

Liu, M., Ene, T.-D., Kirby, R., Cheng, C., Pinckney, N., Liang, R., Alben, J., Anand, H., Banerjee, S., Bayraktaroglu, I., et al. Chipnemo: Domain-adapted llms for chip design. *arXiv preprint arXiv:2311.00176*, 2023.

Liu, S., Fang, W., Lu, Y., Wang, J., Zhang, Q., Zhang, H., and Xie, Z. RTLCoder: Fully open-source and efficient LLM-assisted RTL code generation technique. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2024a.

Liu, S., Lu, Y., Fang, W., Li, M., and Xie, Z. Openllm-rtl: Open dataset and benchmark for llm-aided design rtl generation. In *2024 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pp. 1–9. IEEE, 2024b.

Lu, Y., Liu, S., Zhang, Q., and Xie, Z. Rtllm: An open-source benchmark for design rtl generation with large language model. In *29th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 722–727. IEEE, 2024.

Mi, Z., Zheng, R., Zhong, H., Sun, Y., Kneeland, S., Moitra, S., Kutzer, K., Xu, Z., and Huang, S. Coopetitivev: Leveraging llm-powered coopetitive multi-agent prompting for high-quality verilog generation. *arXiv preprint arXiv:2412.11014*, 2024.

Pinckney, N., Batten, C., Liu, M., Ren, H., and Khailany, B. Revisiting verilogeval: A year of improvements in large-language models for hardware code generation. *ACM Transactions on Design Automation of Electronic Systems*, 2025.

Ping, H., Li, S., Zhang, P., Cheng, A., Duan, S., Kanakaris, N., Xiao, X., Yang, W., Nazarian, S., Irimia, A., and Bogdan, P. Hdlcore: A training-free framework for mitigating hallucinations in llm-generated hdl. *arXiv preprint arXiv:2503.16528*, 2025.

Sun, W., Li, B., Zhang, G. L., Yin, X., Zhuo, C., and Schlichtmann, U. Paradigm-based automatic HDL code generation using LLMs. *arXiv preprint arXiv:2501.12702*, 2025.

Teng, F., Pan, M., Zhang, X., He, Z., Yang, Y., Chai, X., Qi, M., Lu, L., and Yin, J. Verirl: Boosting the llm-based verilog code generation via reinforcement learning. *arXiv preprint arXiv:2508.18462*, 2025.

Thakur, S., Blocklove, J., Pearce, H., Tan, B., Garg, S., and Karri, R. Autochip: Automating hdl generation using llm feedback. *arXiv preprint arXiv:2311.04887*, 2023.

Wang, J., Wang, J., Athiwaratkun, B., Zhang, C., and Zou, J. Mixture-of-agents enhances large language model capabilities. *arXiv preprint arXiv:2406.04692*, 2024a.

Wang, N., Yao, B., Zhou, J., Wang, X., Jiang, Z., and Guan, N. Large language model for Verilog generation with code-structure-guided reinforcement learning. *arXiv preprint arXiv:2407.18271*, 2024b.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35: 24824–24837, 2022.

Xiong, C., Huang, H., Zhao, R., and Chen, D. Hlspi-lot: Llm-based high-level synthesis. *arXiv preprint arXiv:2408.06810*, 2024.

Xu, K., Zhang, G. L., Yin, X., Zhuo, C., Schlichtmann, U., and Li, B. Automated c/c++ program repair for high-level synthesis via large language models. *arXiv preprint arXiv:2407.03889*, 2024.

Yang, Y., Teng, F., Liu, P., Qi, M., Lv, C., Li, J., Zhang, X., and He, Z. Haven: Hallucination-mitigated llm for verilog code generation aligned with hdl engineers. *arXiv preprint arXiv:2501.04908*, 2025.

Zhao, Y., Zhang, H., Huang, H., Yu, Z., and Zhao, J. MAGE: A multi-agent engine for automated RTL code generation. *arXiv preprint arXiv:2412.07822*, 2024.