



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

Heng Ping<sup>1</sup>, Arijit Bhattacharjee<sup>2</sup>, Peiyu Zhang<sup>1</sup>, Shixuan Li<sup>1</sup>, Wei Yang<sup>1</sup>, Anzhe Cheng<sup>1</sup>, Xiaole Zhang<sup>1</sup>, Jesse Thomason<sup>1</sup>, Ali Jannesari<sup>2</sup>, Nesreen Ahmed<sup>3</sup>, Paul Bogdan<sup>1</sup>

<sup>1</sup>University of Southern California



<sup>2</sup>Iowa State University



<sup>3</sup>Cisco AI Research



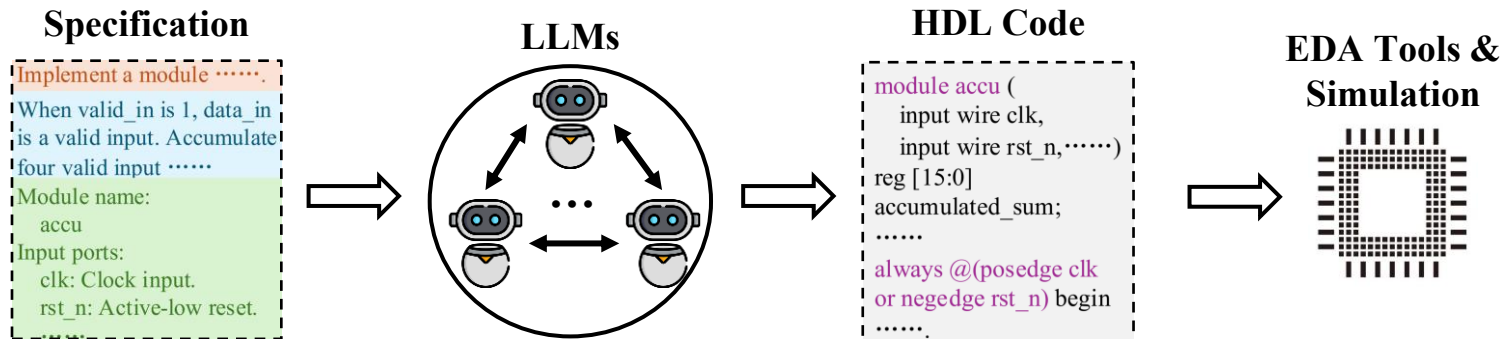
*Proceedings of the 9th MLSys Conference, Bellevue, WA, USA, 2026*

# Background & Motivation



## What is the task?

- Automatically generate HDL (Hardware Design Language) code from natural-language specifications.



## Why HDL is hard for LLMs?



**Sparse Pretraining Data**  
HDL  $\ll$  C++/Python in LLM pretraining corpora






**Concurrent/Timing Semantics**  
Not sequential like common programming language



**Synthesis Constraints**  
HDL must be synthesizable, not just functionally correct

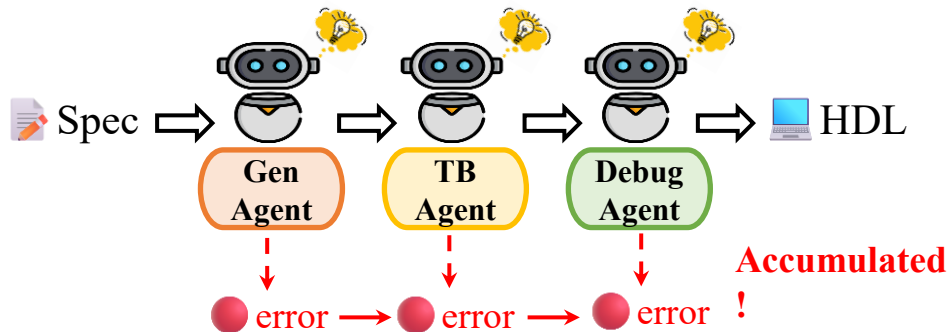
- **General-purpose LLMs lack enough domain-specific knowledge to generate correct HDL.**

## Table: Comparison of Existing RTL Generation Approaches

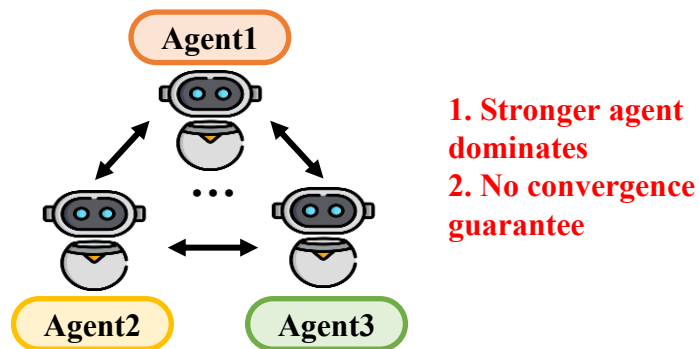
Approach	Representative Works	Limitation
 <b>Prompt Engineering</b>	ParaHDL, AoT, HDLCoRe	✗ Bounded by LLM's sparse HDL knowledge
 <b>SFT / RL</b>	RTLCoder, VeriRL, CodeV	✗ Costly training & data curation
 <b>Multi-Agent Systems</b>	MAGE, CoopetitiveV, VeriMaAS	⚠ Error propagation or chaotic exploration

## Existing Multi-Agent Approaches

### ➤ Sequential Pipeline



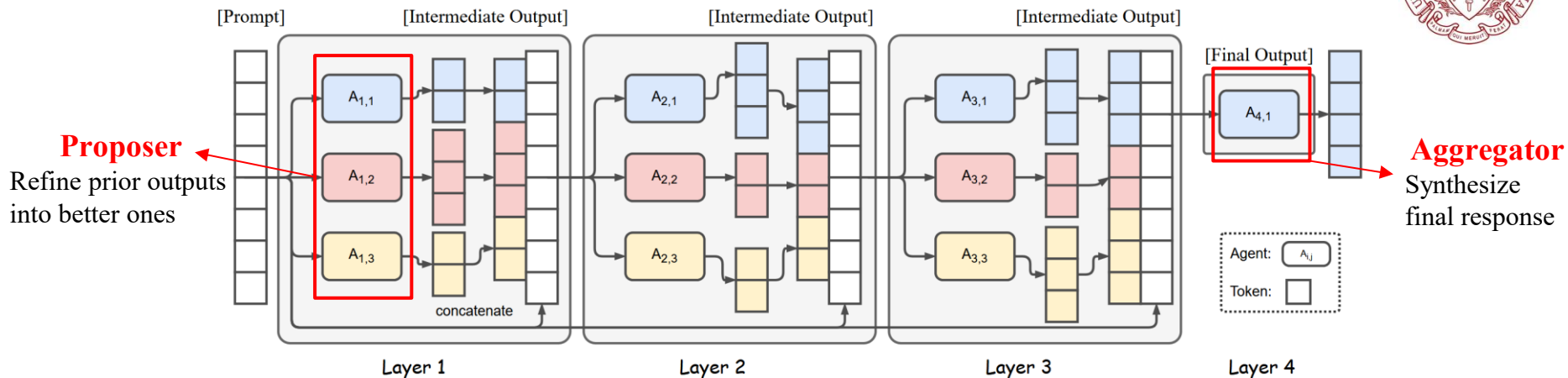
### ➤ Debate Pipeline



# Mixture-of-Agents (MoA)



## MoA Structure — A Layer-wise Paradigm



## Why is MoA promising for HDL?

- ✓ **Layered + parallel** — structured information flow
- ✓ **Superior** on instruction-following & reasoning benchmarks
- 🌟 **Training-free** — works with any LLM backbone

## What does standard MoA still lack?

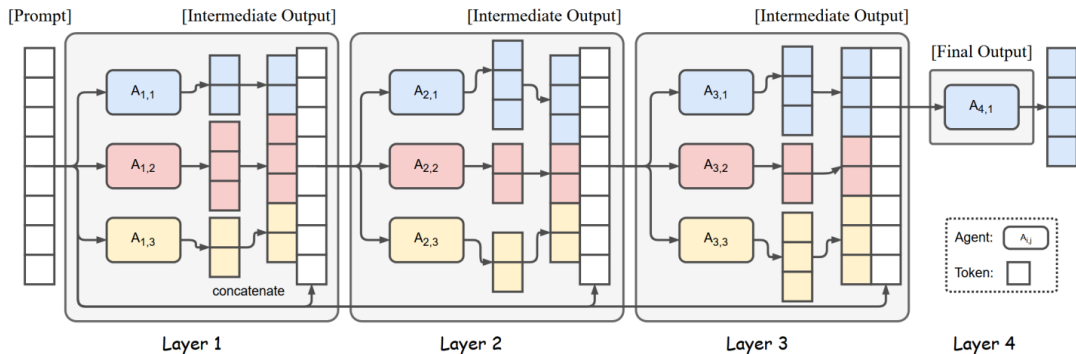
- ✗ **Cascaded dependency**: Each layer only sees the immediately preceding layer
- ✗ **Information loss**: Valuable solutions from earlier layers are forgotten
- ✗ **Limited reasoning path** — all proposers reason similarly, limiting solution diversity

# Key Insight — What Drives MoA Performance?



## Principle of MoA

- MoA performance is highly influenced by the quality and diversity of the agents in MoA layers.<sup>1</sup>



$$\text{Performance} \approx \alpha \cdot \text{Quality} + \beta \cdot \text{Diversity} + \gamma$$

$\alpha > \beta$  (quality matters more than diversity, but both are essential)

## How to improve MoA for HDL?

### How to improve quality?

★ **Quality-Guided Caching**  
Cache and rank all intermediate outputs across layers.

### How to increase diversity?

🌐 **Multi-Path Generation**  
Generate via C++ and Python intermediate representations

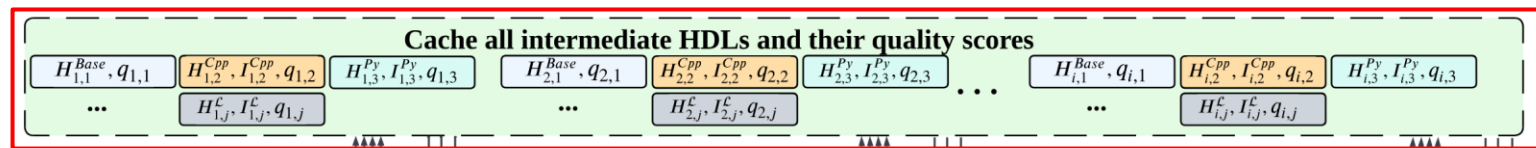
<sup>1</sup>Ref: Rethinking Mixture-of-Agents: Is Mixing Different Large Language Models Beneficial? (2025)

# Framework Overview



## Our Proposed VeriMoA

Global cache indexes ALL intermediates by quality

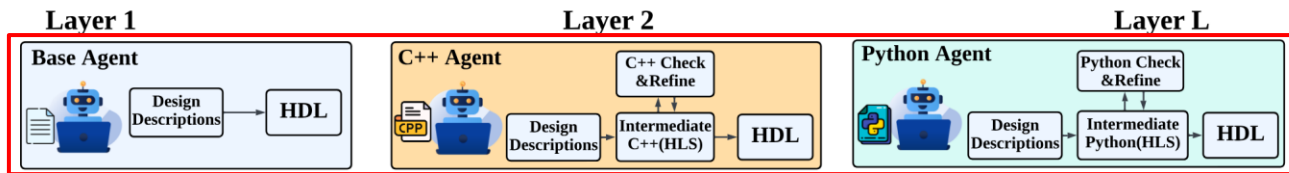
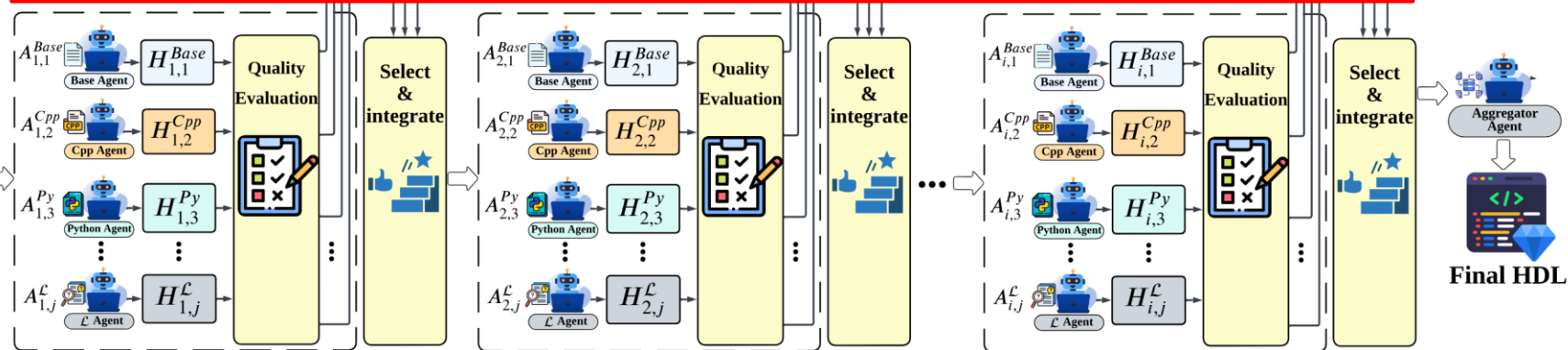


Implement a module of an 8-bit adder .....

The module utilizes a series of bit-level adders .....

Module name: adder\_8bit

Input ports: a[7:0]: 8-bit input operand A. ....



3 agent types to choose per layer

- Three agent types (Base, C++, Python) generate HDL in parallel across layers
- Intermediate outputs are scored, cached globally, and the top-N are passed to deeper layers 6

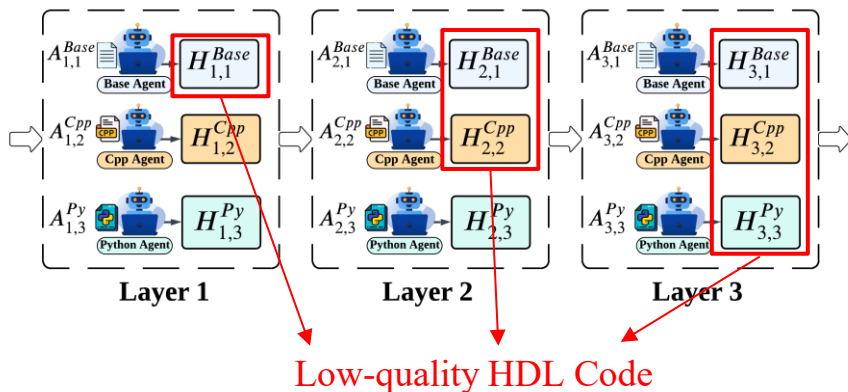
# Innovation 1 — Quality-Guided Caching



## Quality-Guided Caching Mechanism

- Cache all the intermediate HDLs and evaluate the quality of each HDL
- Select top HDLs by quality from cache and integrate them into the prompt for the following agents

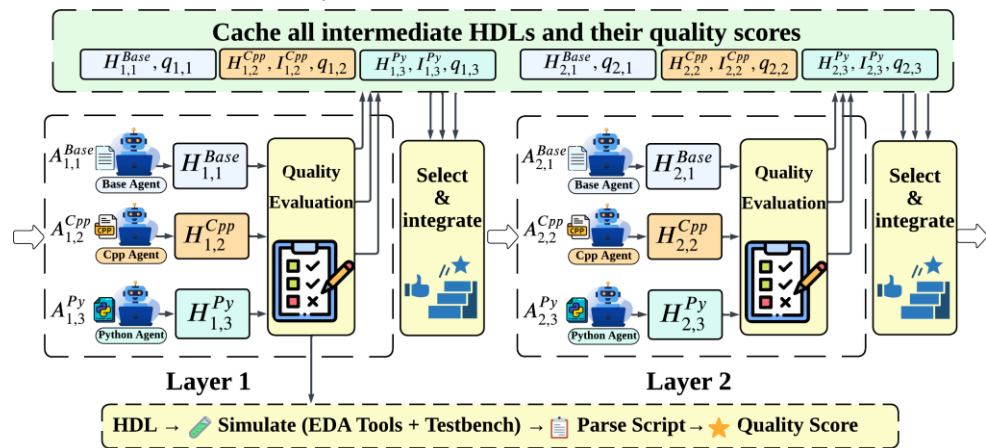
### Standard MoA



#### ⚠ Limitations

- Errors cascade across layers
- Good solutions get lost
- Deeper agents suffer more

### Quality-Guided Cached MoA



#### ✅ Benefits

- Errors filtered at every layer
- Top-quality solutions preserved
- Deeper agents build on the best

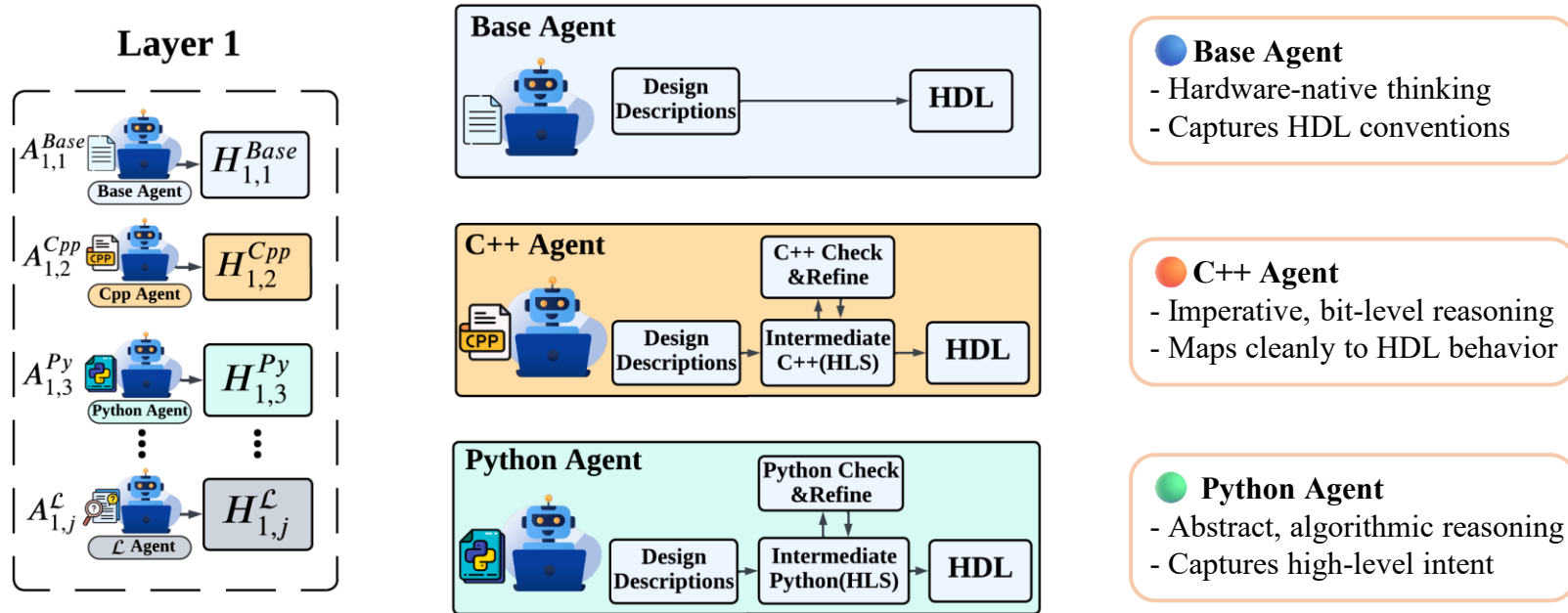
- Global caching + quality scoring → monotonic quality improvement across layers.

# Innovation 2 — Multi-Path Generation



## Multi-path Two-stage Generation

- A multi-path, two-stage approach uses C++/Python as intermediate representations to exploit LLM strengths in high-resource languages while promoting solution diversity



- **Why this works:** LLMs are **far more fluent in C++/Python** → leverage that strength
- **Diversity guarantee:** **Heterogeneous reasoning trajectories** → broader solution space exploration

# Experimental Setup



## Benchmarks & Baselines



### Benchmarks

- **VerilogEval 2.0:**  
156 problems
- **RTL2.0:**  
50 complex designs



### Baselines

- **Prompting LLMs**
- **SFT / RL-based Methods**
- **Multi-agent LLMs**

## LLM Backbones



### LLM Models

- **Open-Source:**
  - Qwen2.5 (7B/14B/32B)
  - Qwen2.5-Coder (7B/14B/32B)
- **Closed-Source:**
  - GPT-4o-mini
  - GPT-4o
- **Open-Source Domain-Specific:**
  - VeriRL-CodeQwen2.5

## Configuration



### Para. & Testing

- **L = 4 layers**
- **M = 6 agents/layer** (2 Base + 2 C++ + 2 Python)
- **Simulator:** Icarus Verilog
- **Metrics:** Pass@1, Pass@3, Pass@5

## Research Questions

**RQ1: Main performance?**  
**RQ2: Contribution of core components?**

**RQ3: Inference cost comparison?**  
**RQ4: Parameter sensitivity?**  
**RQ5: Robustness without manual testbenches?**

# Main Result — vs Non-Training Baselines (RQ1)



## Comparison with Agentic Baselines

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 $\uparrow 3.97$	35.33 $\uparrow 3.13$	42.58 $\uparrow 2.19$	27.49 $\uparrow 8.50$	34.76 $\uparrow 4.23$	39.95 $\uparrow 3.44$
	HDLCoRe	30.92 $\uparrow 8.02$	37.86 $\uparrow 5.66$	46.49 $\uparrow 6.10$	34.68 $\uparrow 15.69$	40.32 $\uparrow 9.79$	44.82 $\uparrow 8.31$
	VeriMaAS	32.81 $\uparrow 9.91$	40.25 $\uparrow 8.05$	49.07 $\uparrow 8.68$	40.33 $\uparrow 21.34$	47.24 $\uparrow 16.71$	50.47 $\uparrow 13.96$
	<b>VERIMO A</b>	<b>56.44 <math>\uparrow 33.54</math></b>	<b>62.07 <math>\uparrow 29.87</math></b>	<b>65.27 <math>\uparrow 24.88</math></b>	<b>52.07 <math>\uparrow 33.08</math></b>	<b>58.10 <math>\uparrow 27.57</math></b>	<b>60.46 <math>\uparrow 23.95</math></b>
Qwen2.5-Coder-7B	Direct	33.91	43.90	48.06	26.16	36.94	41.11
	CoT	36.74 $\uparrow 2.83$	45.88 $\uparrow 1.98$	50.99 $\uparrow 2.93$	31.87 $\uparrow 5.71$	39.36 $\uparrow 2.42$	45.36 $\uparrow 4.25$
	HDLCoRe	40.67 $\uparrow 6.76$	47.81 $\uparrow 3.91$	52.64 $\uparrow 4.58$	36.26 $\uparrow 10.10$	45.61 $\uparrow 8.67$	51.12 $\uparrow 10.01$
	VeriMaAS	44.73 $\uparrow 10.82$	51.62 $\uparrow 7.72$	53.78 $\uparrow 7.72$	44.71 $\uparrow 18.55$	51.37 $\uparrow 14.43$	56.46 $\uparrow 15.35$
	<b>VERIMO A</b>	<b>60.96 <math>\uparrow 27.05</math></b>	<b>68.13 <math>\uparrow 24.23</math></b>	<b>70.69 <math>\uparrow 22.63</math></b>	<b>54.43 <math>\uparrow 28.27</math></b>	<b>62.52 <math>\uparrow 25.58</math></b>	<b>66.24 <math>\uparrow 25.13</math></b>
Qwen2.5-14B	Direct	39.26	49.98	53.74	31.71	40.20	45.87
	CoT	41.24 $\uparrow 1.98$	51.62 $\uparrow 1.64$	54.37 $\uparrow 0.63$	36.49 $\uparrow 4.78$	44.64 $\uparrow 4.44$	48.29 $\uparrow 2.42$
	HDLCoRe	45.18 $\uparrow 5.92$	55.83 $\uparrow 5.85$	59.42 $\uparrow 5.68$	41.78 $\uparrow 10.07$	49.97 $\uparrow 9.77$	53.17 $\uparrow 7.30$
	VeriMaAS	48.97 $\uparrow 9.71$	57.84 $\uparrow 7.86$	61.77 $\uparrow 8.03$	48.14 $\uparrow 16.43$	54.28 $\uparrow 14.08$	58.76 $\uparrow 12.89$
	<b>VERIMO A</b>	<b>65.27 <math>\uparrow 26.01</math></b>	<b>71.61 <math>\uparrow 21.63</math></b>	<b>74.64 <math>\uparrow 20.90</math></b>	<b>55.06 <math>\uparrow 23.35</math></b>	<b>62.76 <math>\uparrow 22.56</math></b>	<b>66.33 <math>\uparrow 20.46</math></b>
Qwen2.5-Coder-14B	Direct	39.74	49.72	52.21	35.61	42.34	46.43
	CoT	43.81 $\uparrow 4.07$	51.27 $\uparrow 1.55$	54.11 $\uparrow 1.90$	40.65 $\uparrow 5.04$	48.77 $\uparrow 6.43$	51.72 $\uparrow 5.29$
	HDLCoRe	46.82 $\uparrow 7.08$	54.79 $\uparrow 5.07$	57.04 $\uparrow 4.83$	47.86 $\uparrow 12.25$	55.19 $\uparrow 12.85$	58.18 $\uparrow 11.75$
	VeriMaAS	50.96 $\uparrow 11.22$	58.48 $\uparrow 8.76$	62.63 $\uparrow 10.42$	51.20 $\uparrow 15.59$	58.93 $\uparrow 16.59$	61.67 $\uparrow 15.24$
	<b>VERIMO A</b>	<b>66.86 <math>\uparrow 27.12</math></b>	<b>73.24 <math>\uparrow 23.52</math></b>	<b>76.38 <math>\uparrow 24.17</math></b>	<b>59.76 <math>\uparrow 24.15</math></b>	<b>64.33 <math>\uparrow 21.99</math></b>	<b>67.22 <math>\uparrow 20.79</math></b>
Qwen2.5-32B	Direct	46.85	56.11	58.41	43.62	48.73	50.75
	CoT	48.72 $\uparrow 1.87$	57.25 $\uparrow 1.14$	59.80 $\uparrow 1.39$	46.99 $\uparrow 3.37$	50.50 $\uparrow 1.77$	52.46 $\uparrow 1.71$
	HDLCoRe	51.75 $\uparrow 4.90$	59.63 $\uparrow 3.52$	61.79 $\uparrow 3.38$	50.95 $\uparrow 7.33$	54.06 $\uparrow 5.33$	58.98 $\uparrow 8.23$
	VeriMaAS	53.57 $\uparrow 6.72$	61.82 $\uparrow 5.71$	64.26 $\uparrow 5.85$	54.18 $\uparrow 10.56$	57.21 $\uparrow 8.48$	61.54 $\uparrow 10.79$
	<b>VERIMO A</b>	<b>71.85 <math>\uparrow 25.00</math></b>	<b>76.48 <math>\uparrow 20.37</math></b>	<b>78.16 <math>\uparrow 19.75</math></b>	<b>63.81 <math>\uparrow 20.19</math></b>	<b>67.60 <math>\uparrow 18.87</math></b>	<b>69.88 <math>\uparrow 19.13</math></b>
Qwen2.5-Coder-32B	Direct	46.93	55.73	57.12	40.40	45.42	48.24
	CoT	48.65 $\uparrow 1.72$	56.31 $\uparrow 0.58$	59.34 $\uparrow 2.22$	45.72 $\uparrow 5.32$	48.75 $\uparrow 3.33$	50.29 $\uparrow 2.05$
	HDLCoRe	51.28 $\uparrow 4.35$	58.63 $\uparrow 2.90$	61.47 $\uparrow 4.35$	51.72 $\uparrow 11.32$	54.73 $\uparrow 9.31$	59.64 $\uparrow 11.40$
	VeriMaAS	56.67 $\uparrow 9.74$	63.46 $\uparrow 7.73$	66.92 $\uparrow 9.80$	55.82 $\uparrow 15.42$	59.97 $\uparrow 14.55$	62.31 $\uparrow 14.07$
	<b>VERIMO A</b>	<b>73.31 <math>\uparrow 26.38</math></b>	<b>79.05 <math>\uparrow 23.32</math></b>	<b>81.27 <math>\uparrow 24.15</math></b>	<b>65.49 <math>\uparrow 25.09</math></b>	<b>71.42 <math>\uparrow 26.00</math></b>	<b>74.11 <math>\uparrow 25.87</math></b>
GPT-4o-mini	Direct	48.97	56.94	58.62	46.60	50.71	51.95
	CoT	52.20 $\uparrow 3.23$	59.83 $\uparrow 2.89$	60.93 $\uparrow 2.31$	48.27 $\uparrow 1.67$	53.56 $\uparrow 2.85$	55.48 $\uparrow 3.53$
	HDLCoRe	53.50 $\uparrow 4.53$	61.46 $\uparrow 4.52$	63.72 $\uparrow 5.10$	51.19 $\uparrow 4.59$	56.72 $\uparrow 6.01$	59.01 $\uparrow 7.06$
	VeriMaAS	57.24 $\uparrow 8.27$	64.77 $\uparrow 7.83$	66.85 $\uparrow 8.23$	57.25 $\uparrow 10.65$	61.46 $\uparrow 10.75$	63.17 $\uparrow 11.22$
	<b>VERIMO A</b>	<b>72.43 <math>\uparrow 23.46</math></b>	<b>76.94 <math>\uparrow 20.00</math></b>	<b>79.46 <math>\uparrow 20.84</math></b>	<b>64.23 <math>\uparrow 17.63</math></b>	<b>67.45 <math>\uparrow 16.74</math></b>	<b>68.67 <math>\uparrow 16.72</math></b>
GPT-4o	Direct	64.74	69.89	71.66	52.48	56.18	57.62
	CoT	66.38 $\uparrow 1.64$	71.18 $\uparrow 1.29$	74.22 $\uparrow 2.56$	54.89 $\uparrow 2.41$	59.22 $\uparrow 3.04$	61.77 $\uparrow 4.15$
	HDLCoRe	69.60 $\uparrow 4.86$	73.52 $\uparrow 3.64$	75.86 $\uparrow 4.20$	56.27 $\uparrow 3.79$	61.47 $\uparrow 5.29$	63.23 $\uparrow 5.61$
	VeriMaAS	71.34 $\uparrow 6.60$	76.12 $\uparrow 6.23$	79.21 $\uparrow 7.55$	61.70 $\uparrow 9.22$	67.25 $\uparrow 11.07$	68.84 $\uparrow 11.22$
	<b>VERIMO A</b>	<b>84.97 <math>\uparrow 20.23</math></b>	<b>89.65 <math>\uparrow 19.76</math></b>	<b>91.03 <math>\uparrow 19.37</math></b>	<b>69.17 <math>\uparrow 16.69</math></b>	<b>73.68 <math>\uparrow 17.50</math></b>	<b>75.62 <math>\uparrow 18.00</math></b>

### Key Findings

- ① +15 to +33 Pass@1 points improvement across all backbones
- ② Smaller models beat larger baselines: VeriMoA-Qwen7B (56.4%) > VeriMaAS-Qwen32B (53.6%)
- ③ Bigger boost for smaller models: Qwen-7B gains +33 pts vs GPT-4o gains +20 pts

**Architectural design beats scale — every backbone improves, smaller ones improve more.**

# Main Result — vs SFT/RL Models + Ablation (RQ1 & RQ2)

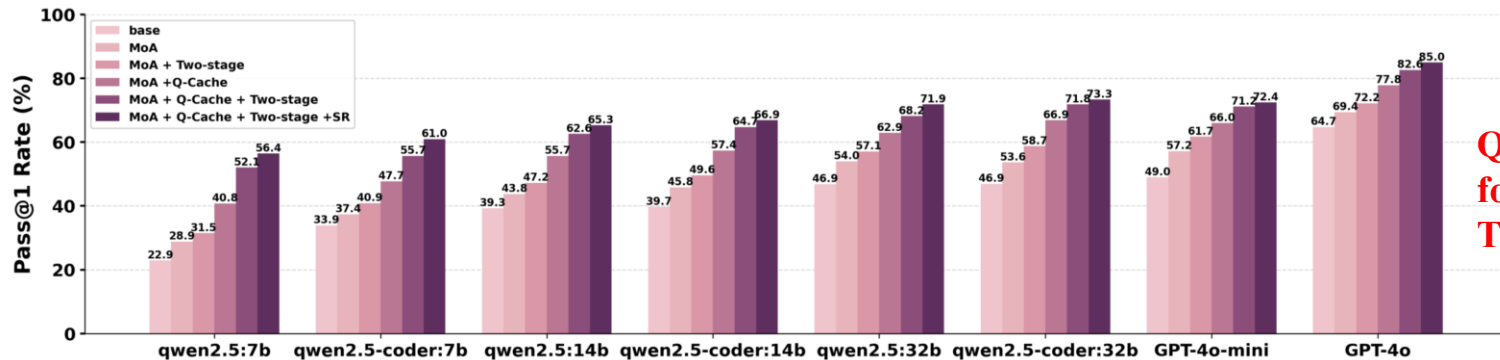


## Comparison with SFT/RL Baselines

	Model	Size	VerilogEval 2.0 (%)			RTL2.0 (%)		
			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
	<b>VeriRL-CodeQwen2.5</b>	7B	<b>66.28</b>	<b>71.35</b>	<b>73.43</b>	<b>61.53</b>	<b>65.27</b>	<b>68.94</b>
Ours	Qwen2.5-Coder-7B + VERIMO A	7B	60.96	68.13	70.69	54.43	62.52	66.24
	Qwen2.5-Coder-14B + VERIMO A	14B	66.86	73.24	76.38	59.76	64.33	66.22
	Qwen2.5-Coder-32B + VERIMO A	32B	<b>73.31</b>	<b>79.05</b>	<b>81.27</b>	<b>65.49</b>	<b>71.42</b>	<b>74.11</b>
	<b>VeriRL-CodeQwen2.5 + VERIMO A</b>	7B	<b>82.47</b>	<b>87.91</b>	<b>90.68</b>	<b>74.45</b>	<b>78.63</b>	<b>82.02</b>

➤ **Training-free framework matches or exceeds fine-tuned models AND is complementary to them.**

## Contribution of Core Mechanisms



**Q-Cache is the foundation unlocking Two-stage's potential.**

# Inference Cost Analysis (RQ3)



## Inference Cost Comparison

Dataset	Model	Baseline		VeriMaAS		VERIMOA-Lite		VERIMOA	
		Pass@1	Tokens	Pass@1	Tokens	Pass@1	Tokens	Pass@1	Tokens
VerilogEval 2.0	Qwen2.5-Coder-32B	46.93%	0.63k	56.67%	3.67k (5.83×)	67.93%	3.73k (5.92×)	73.31%	6.95k (11.03×)
	GPT-4o-mini	48.97%	0.52k	57.24%	3.18k (6.12×)	68.41%	2.94k (5.65×)	72.43%	5.51k (10.60×)
RTLLM 2.0	Qwen2.5-Coder-32B	40.40%	0.76k	55.82%	4.59k (6.04×)	62.82%	4.05k (5.33×)	65.49%	7.28k (9.58×)
	GPT-4o-mini	46.60%	0.67k	57.25%	4.32k (6.45×)	61.29%	3.79k (5.66×)	64.23%	6.47k (9.66×)

⚡ **Cost:** ~10× baseline tokens →  
**Gain:** +20–25 Pass@1 points

🍂 **VeriMoA-Lite:** same cost as  
VeriMaAS, +11 points higher Pass@1

## Accuracy Vs. Tokens——Scale up

➤ varying agents per layer using GPT-4o-mini

Agents	VerilogEval 2.0		RTLLM 2.0	
	Pass@1	Tokens	Pass@1	Tokens
2	64.52%	1.86k (3.57×)	57.89%	2.28k (3.40×)
3	68.41%	2.94k (5.65×)	61.29%	3.79k (5.66×)
4	70.28%	3.88k (7.46×)	62.76%	4.71k (7.03×)
5	71.49%	4.70k (9.04×)	63.58%	5.60k (8.36×)
6	72.43%	5.51k (10.60×)	64.23%	6.47k (9.66×)

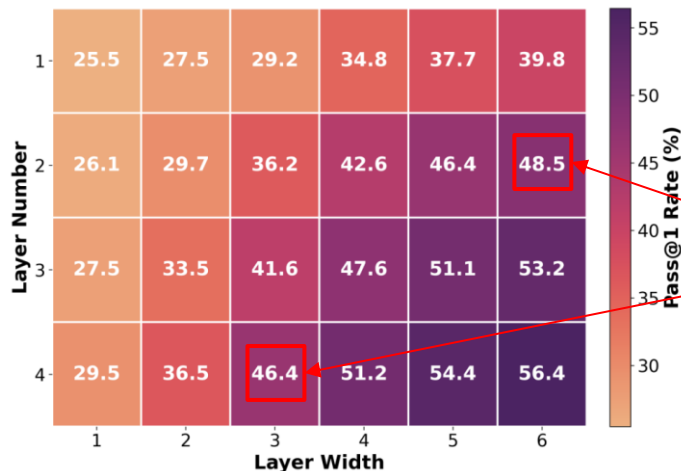
📈 **Cost scales linearly, gains saturate:**  
tokens grow linearly, Pass@1 gains shrink  
from +3.9 to +0.9 pts

💡 **Diminishing returns** — fewer agents  
already capture most of the gain

# Parameter Sensitivity & Robustness (RQ4 & RQ5)



Parameter Sensitivity—*Layer Depth* × *Layer Width* (VerilogEval 2.0 using Qwen2.5:7b)



**Both depth & width matter** —  
1-layer or 1-agent configs cap at <40% Pass@1  
 **Width > Depth at equal budget** —  
2L×6W (48.5%) beats 4L×3W (46.4%)

## Manual Testbench vs LLM-Generated Testbench (GPT-4o-mini)




Testbench	VerilogEval 2.0 (%)			RTL2.0 (%)		
	P@1	P@3	P@5	P@1	P@3	P@5
Golden	72.43	76.94	79.46	64.23	67.45	68.67
LLM-Gen.	67.84	73.57	76.63	60.51	63.87	65.74
Δ	-4.59	-3.37	-2.83	-3.72	-3.58	-2.93

**Minor degradation: only -2.8 to -4.6 Pass@1 pts without manual TB**  
 **Still beats VeriMaAS (which uses manual TB)**



# Conclusion & Future Plan



## Conclusion

-  **Quality-Guided Caching:** monotonic knowledge accumulation across layers
-  **Multi-Path Generation:** leverage LLMs' C++/Python fluency + structured diversity
-  **15–30% Pass@1 gains** across diverse LLMs; small models match large; complementary to fine-tuning

## Future Plan

-  **Hierarchical Generation:** Planner agent decomposes large designs (RISC-V cores, NoC routers) into modules for VeriMoA
-  **Cross-Domain Transfer:** Apply quality-guided multi-path MoA to other code generation tasks (CUDA, assembly, HLS)



**Thanks for Your Listening !**

**Q&A**