

MeltRTL: Multi-Expert LLMs with Inference-time Intervention for RTL Code Generation

Nowfel Mashnoor, Mohammad Akyash, Hadi Kamali, Kimia Azar

Department of Electrical and Computer Engineering (ECE), University of Central Florida, Orlando, FL 32816, USA
 {nowfel.mashnoor, mohammad.akyash, kamali, azar}@ucf.edu

Abstract—The automated generation of hardware register-transfer level (RTL) code with large language models (LLMs) shows promise, yet current solutions struggle to produce syntactically and functionally correct code for complex digital designs. This paper introduces *MeltRTL*, a novel framework that integrates multi-expert attention with inference-time intervention (ITI) to significantly improve LLM-based RTL code generation accuracy without retraining the base model. *MeltRTL* introduces three key innovations: (1) A multi-expert attention architecture that dynamically routes design specifications to specialized expert networks, enabling targeted reasoning across various hardware categories; (2) An inference-time intervention mechanism that employs non-linear probes to detect and correct hardware-specific inaccuracies during generation; and (3) An efficient intervention framework that selectively operates on expert-specific attention heads with minimal computational overhead. We evaluate *MeltRTL* on the VerilogEval benchmark, achieving 96% synthesizability and 60% functional correctness, compared to the base LLM’s 85.3% and 45.3%, respectively. These improvements are obtained entirely at inference time, with only 27% computational overhead and no model fine-tuning, making *MeltRTL* immediately deployable on existing pre-trained LLMs. Ablation studies further show the complementary benefits of multi-expert architecture and ITI, highlighting their synergistic effects when combined.¹

Index Terms—Large language models, RTL code generation, probing, inference-time intervention

I. INTRODUCTION

The integration of large language models (LLMs) into hardware design workflows has sparked growing interest in their ability to translate high-level textual specifications into register-transfer level (RTL) implementations [1]–[3]. By shrinking spec.-to-code time, these models promise faster (and more reliable) prototyping (ultimately a shorter semiconductor time-to-market or TTM) [4]. Unlike conventional software generation, however, RTL is a zero-tolerance domain, where designs must precisely capture cycle-accurate behavior while conforming to structural and verification constraints [5]. In practice, current LLM-generated RTL often compiles but deviates from the specification, e.g., reflecting mismatches, incorrect clocking, missing reset/handshakes, etc. These discrepancies, often categorized as *hallucinations* [6], [7], compromise functional correctness, reliability, and even security [8].

Reliability in language generation is commonly improved via grounding models with external knowledge [9], retrieval-augmented generation [10], and tool integration [11]. While effective as factual anchors, such methods leave the model’s internal computation (e.g., encoding and propagation) unchanged.

Recent studies, on the other hand, point to the benefit of directly steering internal representations for truthfulness. Techniques such as activation editing [12], [13] and intermediate-level alignment [14], [15] demonstrate that errors often arise from systematic misalignments within latent spaces, and by reshaping these dynamics at inference, models can be steered toward truthful and instruction-consistent outputs [16].

To the best of our knowledge, no prior work has applied inference-time, activation-level interventions to RTL generation. Extending these advances to hardware is particularly critical: unlike open-domain text, where hallucinations produce factual inaccuracies, hallucinations in RTL manifest as logic that fails simulation or formal checks, degrades power–performance–area (PPA), or undermines security guarantees [6], [17]. Moreover, the scarcity of high-quality (verified) RTL datasets limits the feasibility of large-scale fine-tuning, making these strategies especially attractive. By steering the internal activations of LLMs during generation, we can directly enforce semantic alignment, ensuring RTL implementations remain reliable, efficient, and instruction-consistent without the prohibitive cost of retraining.

While motivating, leveraging the intermediate layers of LLMs for RTL generation introduces two challenges:

- (i) *Targeting Representation Components*: Identifying which elements of the model’s internal representations are most responsible for functional correctness and should be modified.
- (ii) *Robust and Efficient Intervention*: Designing intervention strategies that are both lightweight and robust, ensuring they effectively enhance RTL functional correctness without introducing instability or excessive computational overhead.

To overcome these challenges, we introduce *MeltRTL*: multi-expert LLMs with inference-time intervention for RTL Code Generation, a framework that couples probe-guided multi-expert attention with inference-time steering, enabling correctness-oriented control without retraining base models. Our contributions are threefold:

- (i) *Probe-Guided Component Identification*: We curate a small dataset of 200 samples of instruction–code pairs categorized by hardware specifications (e.g., FSM, arithmetic, etc.). Using this dataset, we train multiple lightweight classifiers on internal activations to identify category-specific attention heads that most strongly predict functional correctness, revealing structure in how LLMs internalize RTL reasoning.
- (ii) *Representation-Level Steering for Correctness*: Leveraging the identified heads, *MeltRTL* applies targeted, non-

¹Code available at: <https://github.com/mashnoor/melt-rtl>

TABLE I: Comparison of Prior LLM-based RTL Generation Models. Legend: ✓ = yes, ✗ = no, ◐ = partial/indirect, — = not reported.

Model	Train free	Activation Steering	Multi Expert	Agentic Tools	RAG / Self-Ref.	Fine Tuning	Open-Source
RTLCoder [1]	✗	✗	✗	✗	✗	✓FT	✓
OriGen [2]	✗	✗	✗	◐	✓	✓LoRA FT	✓
BetterV [18]	✗	✗	✗	✗	◐	✓FT	◐
CodeV [19]	✗	✗	✗	✗	✓	✓FT	◐
CraftRTL [20]	✗	✗	✗	✗	◐	✓FT (synt.)	✓
VerilogCoder [21]	✓	✗	✗	✓	◐	✗	✗
MAGE [22]	✓	✗	✗	✓	◐	✗	✓
RTL++ [23]	✗	✗	✗	✗	✓	✓FT (graph)	◐
HDLCoRe [6]	✓	✗	✗	◐	✓	✗	—
HaVen [7]	✗	✗	✗	✗	✓	✓FT	✓
ScaleRTL [24]	✗	✗	✗	✓	—	✓(reason)	—
VeriSeek [25]	✗	✗	✗	✗	✗	✓RL	—
MeltRTL	✓	✓	✓	◐	◐	✗	✓

linear interventions during decoding to enforce specification-faithful behavior. This reduces RTL-specific hallucinations and improves correctness while remaining computationally lightweight (no base model fine-tuning required).

(iii) Domain-Specialized Multi-Expert Behavior: We show that steering different subsets of heads enables the model to act as an *expert* in specific hardware design domains. This multi-expert specialization equips MeltRTL to adaptively improve syntactic and functional correctness across diverse categories of RTL designs of varying complexity.

II. RELATED WORKS

A. LLM for RTL Code Generation

LLMs show strong capabilities in translating natural language into executable code across software languages, driven by large-scale pretraining on diverse codebases (e.g., Codex [26], CodeGen [27]). Recent efforts extend these advances to hardware, covering design optimization [28], debugging [29], and vulnerability detection [30]. Several studies directly target LLM-based RTL generation [1], [2], [18]–[20], [31]. Early systems [31] rely on prompt engineering and tool feedback (compiler/simulator in the loop) to coax general-purpose LLMs (e.g., GPT-3.5/4) into producing synthesizable code. While helpful, these frameworks are iteration-heavy and often require manual post-processing to ensure semantic fidelity.

To reduce manual intervention, recent research has adapted LLMs specifically for RTL. For instance, RTLCoder [1] augmented scarce datasets by synthesizing instruction–code pairs with GPT-3.5, while OriGen [2] introduced iterative refinement through self-reflection and augmentation. BetterV [18] guided generation toward hardware design objectives such as PPA (power, performance, area). CraftRTL [20] enriched context with auxiliary artifacts like state diagrams and simulation traces, improving structured reasoning. ScaleRTL [24] scaled training to a 3.5B-token reasoning dataset that captures RTL semantics and applied reasoning-oriented techniques for RTL code generation. As summarized in Table I, prior efforts have largely focused on three levers: (i) fine-tuning or leveraging large auxiliary datasets, (ii) orchestrating agentic toolchains with compilers and simulators in the loop, or (iii) enriching prompts and contexts with additional artifacts. In contrast, MeltRTL is, to the best of our knowledge, the first approach to directly intervene in the internal representations of LLMs for

RTL code generation. This activation-level, train-free strategy is complementary to existing methods and can be readily combined with retrieval-augmented or agentic pipelines.

B. Hallucination Mitigation in LLMs

A growing literature studies internal representations to improve truthfulness and reduce hallucinations [32]. While effective in natural language, direct transfer to structured domains like RTL is nontrivial [33], [16], [34].

TruthX [16] employs a contrastive autoencoder to separate factuality from semantics. Although it shows promise in QA, this method requires finely labeled pairs and incurs lossy reconstruction, which complicates adaptation to RTL tasks. CCS [34] finds truth-related directions through unsupervised contrastive consistency, yet hinges on linguistic negation signals absent in RTL codes. ITI [33] ranks attention heads by linear probe accuracy, offering a fast and modular way to attribute truth-related signals. However, its reliance on labeled outputs and the assumption of linear separability limit robustness, and treating heads independently overlooks important joint dependencies. Truth Forest (TrFr) [35] extends ITI with orthogonal probes and lightweight inference-time steering, improving representation coverage. Still, it remains dependent on labeled data and interpretable latent directions, which restricts its applicability in highly structured and domain-specific settings like hardware generation.

Considering these limitations and the promising results of prior work, MeltRTL is designed as a multi-expert framework for RTL code generation. Unlike existing approaches, MeltRTL directly targets correctness-critical attention heads through probe-guided inference-time steering, reducing hallucinations without requiring extensive labeled data or lossy latent representations. As a result, MeltRTL remains both practical and domain-aware, addressing the strict structural and semantic constraints of hardware description languages.

III. PROPOSED METHODOLOGY: MELTRTL

Figure 1 provides an overview of the MeltRTL framework. As shown, it consists of three main stages: (i) constructing and categorizing an instruction–code dataset (small size), (ii) training probes on intermediate activations to identify correctness-critical heads, and (iii) steering chosen heads at inference time to ensure specification-faithful RTL generation.

A. Dataset Collection for Probe Training

To enable effective training of the lightweight probes used in MeltRTL, we curate a small dataset of 200 samples comprising paired design specifications (instructions) and corresponding RTL code implementations. The dataset was collected from a wide range of open-source RTL (.v/.sv) repositories, including GitHub, OpenCores, and other community-maintained RTL repositories (examples spanning arithmetic units, controllers, encoders/decoders, and memory components). A key difference between this dataset preparation and the dataset used for fine-tuning base models is its intentionally relatively small size².

²Probes are lightweight classifiers/regressors trained on model activations; they need clean, labeled signals more than volume.

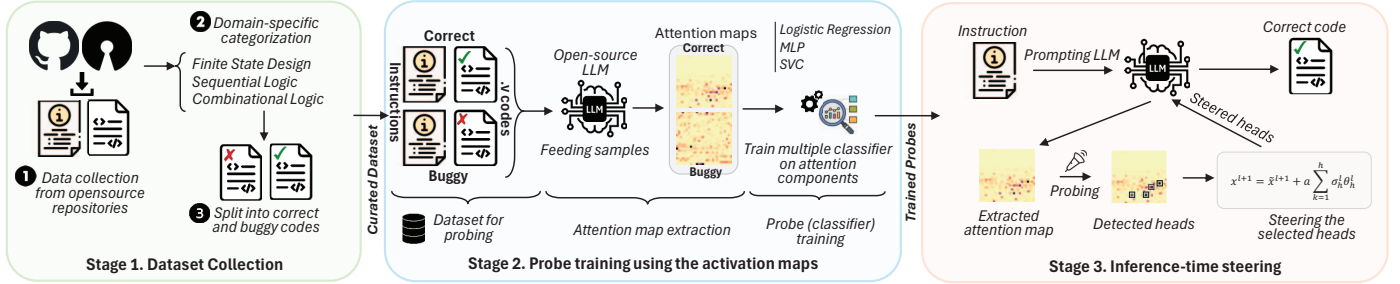


Fig. 1: Overall framework of MeltRTL, consisting of three stages: (i) dataset collection, (ii) detection of heads relevant to functional correctness, and (iii) LLM steering via modification of the identified heads.

From these repositories, we curated a total of **200 representative samples**, each consisting of an instruction–code pair. For every RTL module, a corresponding textual specification or derived instruction was either directly obtained, formulated from documentation and code comments, or generated by an LLM (validated manually for correctness and accuracy). To ensure high-quality supervision, we apply a two-stage filtering pipeline: (i) We applied syntactic validation to remove incomplete or ill-formed modules; and (ii) We employed simulation-based checks and assertion-driven formal verification to assess functional behavior. Based on these evaluations, the dataset is partitioned into two categories:

(i) Functionally Correct Code: RTL codes that compile and satisfy their design intent without observable deviations.

(ii) Functionally Incorrect (Buggy) Code: Modules that are syntactically correct but functionally misaligned, exhibiting errors such as incorrect state transitions, misaligned outputs, or signals with incorrectly initialized values.

Beyond correctness, we further categorize designs into three fundamental logic classes (w.r.t. their application):

(1) Combinational Modules: Stateless modules (w/o register and memory) where outputs depend solely on current inputs.

(2) Sequential (Datapath) Modules: Modules driven by registers or memory elements, especially for building datapath (e.g., datapath in pipeline or systolic arrays).

(3) Finite-State Design (Controller): Modules dominated by explicit state machine transitions or controlling signals.

This additional categorization was motivated by our goal of understanding **which attention heads are critical for which category of design**. By associating probe predictions with functional domains, MeltRTL learns to allocate specialized expertise across different logic types. This categorization directly informed the multi-expert design of our framework, enabling expert-specific intervention strategies that improve generation fidelity across diverse RTL patterns.

B. Classifier (Probe) Training

To train probes that can detect correctness-related signals, we followed a structured pipeline. First, the curated instruction–code pairs were fed into the pre-trained LLM, and we collected the corresponding activation outputs from individual attention heads. From these activations, we isolated the outputs of individual attention heads, where each head provides a representation vector $h_i \in \mathbb{R}^d$. These vectors encode information about intermediate reasoning steps during code generation.

Next, we labeled the head representations according to the functional correctness of the generated code (linking activations to functional outcomes). Head activations originating from modules classified as correct were assigned label $y = 1$, while those from buggy modules were labeled $y = 0$. This transformation provided a training dataset in the form $\{(h_i, y)\}$, suitable for supervised probe training (probes correlate attention-head activity with RTL correctness).

1) Probe Architectures: To capture discriminative signals across representations, we trained multiple classifiers:

(i) Logistic Regression (LR): A linear baseline that tests whether correctness can be separated with a hyperplane.

(ii) Multi-Layer Perceptron (MLP): A shallow non-linear network that captures higher-order interactions missed by LR.

(iii) Support Vector Classifier (SVC): A (RBF) kernel-based margin classifier, useful for separating correctness signals in curved or irregular regions of the space.

Formally, the objective of probe is to learn a specific mapping ($f : h_i \mapsto \{0, 1\}$), where 0 denotes buggy code and 1 denotes functionally correct code. This formalizes the probe’s task as binary classification over head representations. Now, as shown in Equations 1, 2, 3, the model estimates the probability of correctness based on LR (followed by a sigmoid, which is effective when correctness signals align linearly in feature space), MLP (e.g., ReLU allowing the probe to model interactions across features, capturing more subtle correctness cues), and SVC (carving out complex decision boundaries, making it suitable when correctness signals are not linearly separable.) probes, respectively.

$$P(y = 1 | h_i) = \sigma(W^\top h_i + b). \quad (1)$$

$$P(y = 1 | h_i) = \sigma(W_2 \phi(W_1 h_i + b_1) + b_2). \quad (2)$$

$$f(h_i) = \text{sign} \left(\sum_{j=1}^m \alpha_j y_j K(h_i, h_j) + b \right). \quad (3)$$

2) Ensemble and Majority Voting: Each classifier was trained independently, and their predictions were combined into an ensemble using majority voting (Equation 4):

$$\hat{y} = \text{mode}\{f_{\text{LR}}(h_i), f_{\text{MLP}}(h_i), f_{\text{SVC}}(h_i)\}. \quad (4)$$

This ensemble improves robustness by pooling the complementary strengths of linear and non-linear models, reducing sensitivity to the weaknesses of any single probe. Beyond

classification, the probes highlight **which attention heads consistently correlate with correctness**. These correctness-critical heads form the foundation of MeltRTL’s inference-time intervention, where steering is applied selectively to maximize alignment with design intent while keeping overhead low.

C. Inference-Time Steering of Selected Heads

With a ranked set of correctness-critical heads, inference-time steering can be accomplished in three stages: (i) background formulation of multi-head attention, (ii) probe-guided selection of correctness-critical heads, and (iii) targeted residual corrections applied only to those heads during generation.

(i) Attention background: In a transformer layer l , each head $h \in \{1, \dots, H\}$ computes z_h^l as follows:

$$z_h^l = \text{Att}_h^l(P_h^l x^l), \quad (5)$$

where $x^l \in \mathbb{R}^d$ is the residual state, P_h^l the input projection, and Att_h^l the attention function. The heads are projected and aggregated to update the residual stream:

$$\tilde{x}^{l+1} = x^l + \sum_{h=1}^H Q_h^l z_h^l, \quad (6)$$

with Q_h^l denoting the output projection. It defines how information from heads is fused into the model’s hidden state.

(ii) Probe-guided head selection: At inference, the ensemble decision from Section III-B determines whether a head is correctness-critical. Specifically, for head h in layer l , the binary gating variable $\sigma_h^l \in \{0, 1\}$ is set by the probe ensemble, and the steering direction $\theta_h^l \in \mathbb{R}^d$ is a unit-normalized correction vector (nudging activations toward the “functionally correct” region of representation space).

(iii) Inference-time steering: During generation, the residual stream is selectively corrected as:

$$x^{l+1} = \tilde{x}^{l+1} + \alpha \sum_{h=1}^H \sigma_h^l \theta_h^l, \quad (7)$$

where $\alpha > 0$ controls intervention strength. This ensures that only correctness-critical heads contribute adjustments, while all other heads remain untouched.

By design, this mechanism ensures (i) **targeted interventions**, where only heads identified as correctness-critical are adjusted, and (ii) **compatibility**, since the additive corrections are lightweight and preserve the model’s overall distribution while nudging it away from error-prone regions.

IV. RESULTS AND EVALUATION

A. Experimental Setup

We ran all experiments on a GPU server with clustered NVIDIA H100 GPUs for large-scale training and evaluation. The stack included PyTorch with HuggingFace Transformers for model operations and probing, iVerilog for RTL simulation, and Yosys for RTL synthesizability analysis. The setup involved 300 runs across multiple parameter dimensions. Probe architectures included logistic regression, MLPs, and SVCs with RBF kernels, implemented in R. Top-K head selection ranged from

TABLE II: Hyperparameter Configuration. Full Grid Search over all Parameters with 297 Independent Experiments.

Hyperparameter	Values Explored
Probe Type	Linear (LR), Multi-Layer Perceptron (MLP), Support Vector (SVC) with RBF Kernel
Top-K Heads	5, 10, 15, 20, 25, 30, 35, 40, 45, 48
Intervention Strength (α)	0.1, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0
Probe Max Iterations	1000
Probe Regularization	0.01

5–48, and intervention strength from 0.1–5.0. All probes used 1000 iterations with 0.01 regularization. Hyperparameter details are summarized in Table II.

B. Overall Performance Comparison

Our results show MeltRTL substantially improves both synthesizability and functional correctness on the VerilogEval Dataset [36]. Table III compares the baseline LLM (QwenCoder2.5-14B), MeltRTL (single-expert), and MeltRTL (multi-expert). The baseline reached 85.33% synthesizability, MeltRTL (single) improved to 93.33%, and MeltRTL (multi) achieved 96%. Functional correctness showed larger gains. The baseline achieved 45.33%, MeltRTL (single) reached 52%, and MeltRTL (multi) 60%. As shown in Table III, statistical analysis shows that the multi-expert variant was significant ($p = 0.015$, $d = 0.28$), confirming robust improvements, while the single-expert showed less reliable trends.

C. Expert-Specific Analysis

To examine MeltRTL’s adaptability, we analyzed performance across three expert classes: combinational logic, sequential (datapath) logic, and FSMs (controllers). Table IV compares functional correctness when various categories are used for attention heads and steering. As shown, combinational logic provides the highest baseline at 68%, improving to 72.7% with MeltRTL using the combinational expert. Sequential (datapath) logic and FSMs remained difficult, with the base model lowest overall. FSMs were most challenging, with only 12.1% baseline correctness, but MeltRTL’s multi-expert method improved it by approximately 10%.

TABLE III: Comparison of Models (Synthesizability and Functional).

Method	Synth.	Func.	Impr. vs. Base	p-val	d
Base	85.33%	45.33%	-	-	-
MeltRTL (Single)	93.33%	52.00%	+6.67%	0.299	0.12
MeltRTL (Multi)	96.00%	60.00%	+14.67%	0.015	0.28

TABLE IV: Functional Correctness % by Expert. Columns correspond to evaluated models, rows to problem types. Base: QwenCoder2.5-14B w/o steering. General: MeltRTL with non-specific (global) steering. Comb./Seq./FSM: MeltRTL using single-expert heads for comb., seq., or FSM designs. Multi-Expert: MeltRTL leveraging heads from all three categories jointly.

Category	Base	General	Comb.	Seq.	FSM	Multi-Expert
Comb.	68.0%	70.7%	72.7%	71.0%	64.0%	79.8%
Seq.	31.0%	33.3%	28.6%	33.3%	26.2%	43.6%
FSM	12.1%	6.1%	6.1%	15.2%	15.3%	21.7%

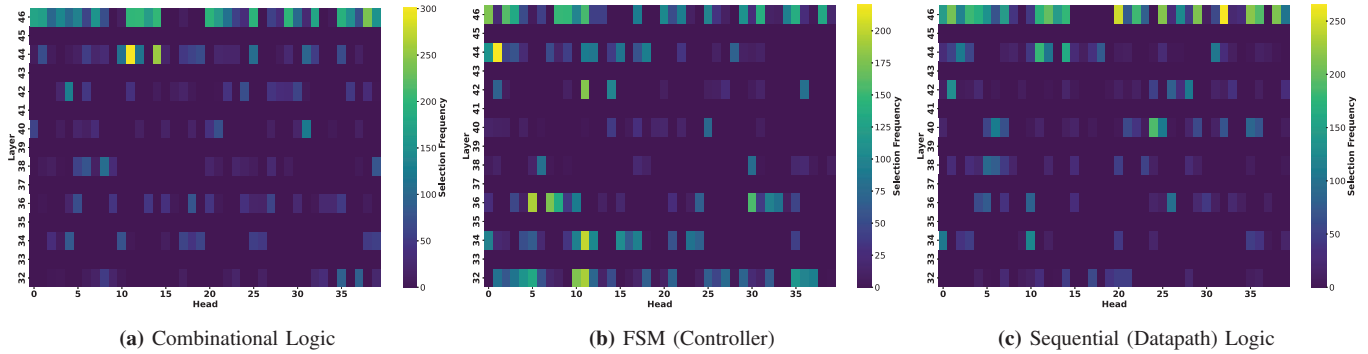


Fig. 2: Heatmaps of Correctness-Critical Attention Heads Selection Across Design Classes. Brighter regions indicate heads more strongly correlated with functional correctness, revealing that distinct subsets of heads specialize by design category

Table IV details these metrics, confirming that domain-matched experts yield consistent gains, while the multi-expert model provides the best results, showcasing the importance of category-based attention head and steering. Improvements scale with design complexity, showing specialized interventions are most beneficial for difficult categories.

D. Probe Performance and Head Selection Analysis

To assess the effectiveness of probe-guided head selection (Section III-B1), we analyzed probe performance and classifier discriminability across RTL design experts. As shown in Table V, we tested three probe architectures. For functional correctness, SVC with RBF consistently led, peaking at $K=15$ and $\alpha=3.0$. LR (linear) probes showed high variance (best 59.33%, worst 33.67%), while Multi-Layer Perceptron (MLP) achieved intermediate, more stable results. Figure 2 reflects domain-specific head activation: combinational logic concentrated in specific clusters, while sequential logic and FSMs showed more distributed patterns, reflecting higher complexity and need for diverse attention. Table VI shows most expert heads (73.7–94.7%) localize in layers 42–46 across probes and domains. This trend is strongest for LR and SVC probes, while MLPs distribute more evenly, suggesting non-linear probes can also capture earlier-layer signals relevant to correctness.

E. Activation Space and Discriminative Analysis

To analyze how probes distinguish correct from incorrect generations, we examined activation patterns across classifier types. Figure 3 shows L2 norm distributions: SVC with

TABLE V: Top and Bottom 3 by Functionality and Synthesizability

Rank	Cls.	K	Str.	Metric	Score
Top 1	SVC with RBF	15	3.0	Func.	60.00
Top 2	SVC with RBF	35	4.5	Func.	59.13
Top 3	LR (Linear)	48	0.1	Func.	59.33
Top 1	Multi-Layer (MLP)	10	2.0	Synth.	99.33
Top 2	LR (Linear)	45	3.5	Synth.	99.33
Top 3	LR (Linear)	5	1.5	Synth.	98.67
Bottom 3	LR (Linear)	15	4.0	Func.	37.67
Bottom 2	LR (Linear)	30	1.0	Func.	35.22
Bottom 1	Multi-Layer (MLP)	30	4.0	Func.	33.67
Bottom 3	Multi-Layer (MLP)	15	4.5	Synth.	81.33
Bottom 2	Multi-Layer (MLP)	48	2.5	Synth.	79.67
Bottom 1	SVC with RBF	45	3.0	Synth.	74.67

TABLE VI: Percentage Distribution of Selected Expert Heads in Mid Layers (32–40) vs. Late Layers (42–46) (Expert heads consistently localize in the late (final) layers, especially L44 and L46).

Model	Expert	Mid (32–40)	Late (42–46)
SVC with RBF	General	11.6	88.4
	Comb.	10.1	89.9
	Seq.	16.9	83.1
	FSM	26.3	73.7
LR (Linear)	General	6.5	93.5
	Comb.	5.3	94.7
	Seq.	7.9	92.1
	FSM	22.8	77.2
Multi-Layer (MLP)	General	53.6	46.4
	Comb.	54.7	45.3
	Seq.	60.8	39.2
	FSM	64.9	35.1

RBF probes give the clearest separation (correct 15–25 vs. incorrect 5–15), LR (linear) probes show overlap but a noticeable shift, and multi-layer (MLP) display multi-peak patterns, indicating nuanced non-linear representations. Figure 4 projects discriminative directions across hardware experts. For combinational logic, correct and incorrect samples separate clearly with distinct centroid positions. Sequential logic shows tighter clusters and smaller margins, matching its lower baseline performance. FSMs exhibit the most overlap, explaining their difficulty even with specialized probes. Overall, these results confirm correctness-related signals in attention activations, beyond probe-specific artifacts.

F. Intervention Optimization Analysis

To optimize MeltRTL’s intervention, we analyzed strength and head selection strategies. Figure 5 shows a non-monotonic relationship between α and performance, peaking at $K=15$ heads. At low strength ($\alpha \leq 1.0$), corrections are too weak for meaningful gains. In the optimal range ($\alpha = 2.5–3.5$), functional correctness improves substantially. Beyond $\alpha = 4.0$, excessive intervention destabilizes generation, as steering signals overwhelm the model’s natural patterns.

G. Computational Overhead Analysis

The computational overhead of ITI is minimal, consisting of a one-time offline setup and a lightweight per-token cost during inference. The setup phase, i.e., activation collection and probe training, is cached and does not affect generation

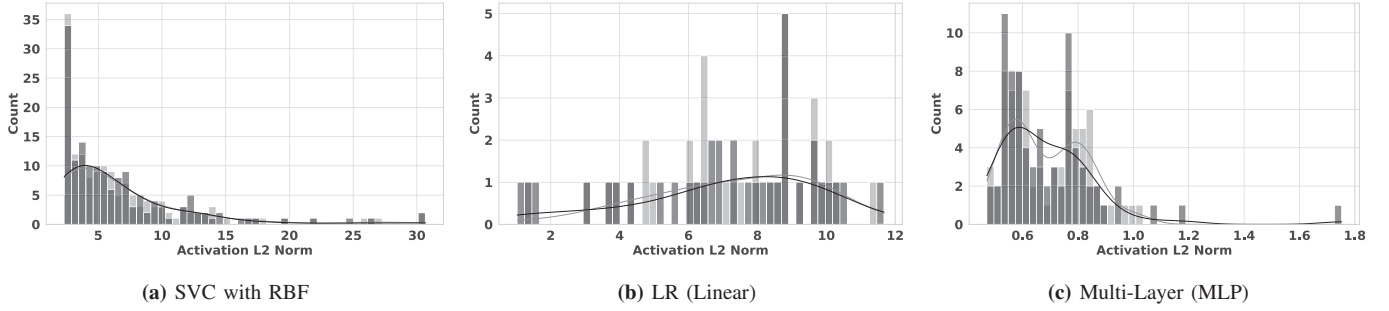


Fig. 3: Distribution of Activation Vector Norms for Correct vs. Incorrect Predictions (for L2 Norm of Attention-Head Activations). The variation in norm ranges highlights how different probe architectures capture correctness-related signals with distinct sensitivity and scaling.

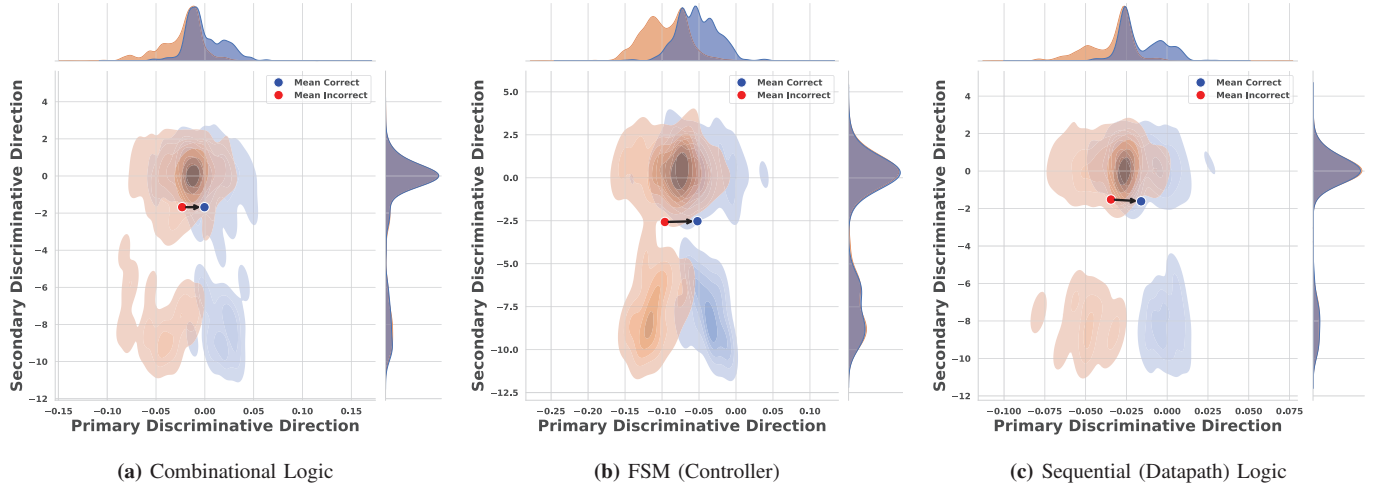


Fig. 4: Activation Shifts Across All Classifiers for Expert Tasks. Combinational designs show the clearest separation with distinct centroids, FSMs exhibit significant overlap, and sequential logic falls in between, reflecting the varying difficulty of correctness alignment across design categories.

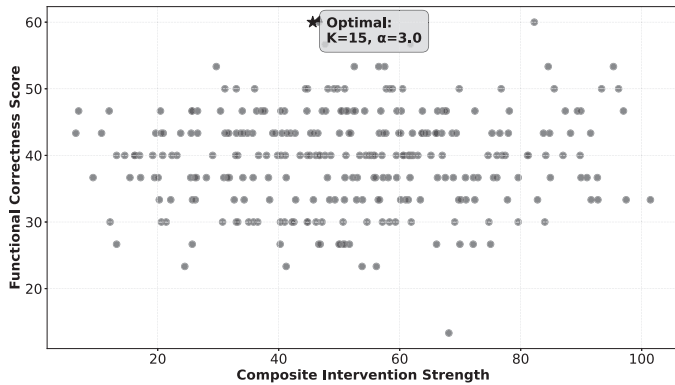


Fig. 5: Intervention strength (α) vs. correctness; peak near $K=15$, $\alpha \approx 3.0$.

latency. At inference, the added costs are selecting the top- k heads (negligible, $O(N_p \log N_p)$) and applying interventions at L_i layers, where each intervention is a vector addition of size d_{model} with cost $O(d_{model})$. Formally, the cost of ITI is

$$C_{ITI} = C_{base} + |L_i| \cdot O(d_{model}), \quad \Delta C = |L_i| \cdot O(d_{model}), \quad (8)$$

where the base model cost is $C_{base} \approx O(s \cdot d_{model}^2)$. Hence, the relative overhead could be calculated as follows:

$$\frac{\Delta C}{C_{base}} \approx \frac{|L_i|}{s \cdot d_{model}}, \quad (9)$$

TABLE VII: Computational Overhead Analysis

Metric	Base	MeltRTL
Avg. Time / sample (s)	11.27	14.27
Std. Dev (s)	5.36	5.83

Equation 9 shows that ITI adds only simple vector additions compared to the dominant quadratic operations of transformers, leading to negligible practical latency. Table VII provides empirical measurements of the timing overhead, showing that MeltRTL adds only 3 seconds per sample (RTL) generation.

V. CONCLUSION

This paper introduced MeltRTL, the first framework to enhance RTL code generation through inference-time interventions at the representation level, without retraining or fine-tuning LLMs. By combining probe-guided head selection with a multi-expert steering mechanism, MeltRTL achieves substantial gains in synthesizability and functional correctness over strong base models (functional correctness from 45.3% to 60.0% while raising synthesizability from 85.3% to 96.0%), with only modest computational overhead. By demonstrating that internal representations of LLMs can be effectively steered at inference time, MeltRTL establishes a new paradigm for LLM-based RTL generation, which (orthogonally) complements existing scaling, fine-tuning, and agentic strategies.

REFERENCES

- [1] S. Liu, W. Fang, Y. Lu, J. Wang, Q. Zhang, H. Zhang, and Z. Xie, "Rtlcoder: Fully open-source and efficient llm-assisted rtl code generation technique," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2024.
- [2] F. Cui, C. Yin, K. Zhou, Y. Xiao, G. Sun, Q. Xu, Q. Guo, Y. Liang, X. Zhang, D. Song *et al.*, "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*, 2024, pp. 1–9.
- [3] M. Z. S. Khan, N. Mashnoor, M. Akyash, K. Azar, and H. Kamali, "Sage-hls: Syntax-aware ast-guided llm for high-level synthesis code generation," 2025. [Online]. Available: <https://arxiv.org/abs/2508.03558>
- [4] H. Huang, Z. Lin, Z. Wang, X. Chen, K. Ding, and J. Zhao, "Towards llm-powered verilog rtl assistant: Self-verification and self-correction," *arXiv preprint arXiv:2406.00115*, 2024.
- [5] M. Akyash, K. Azar, and H. Kamali, "Decortl: A run-time decoding framework for rtl code generation with llms," *arXiv preprint arXiv:2507.02226*, 2025.
- [6] H. Ping, S. Li, P. Zhang, A. Cheng, S. Duan, N. Kanakaris, X. Xiao, W. Yang, S. Nazarian, A. Irimia *et al.*, "Hdlcore: A training-free framework for mitigating hallucinations in llm-generated hdl," *arXiv preprint arXiv:2503.16528*, 2025.
- [7] Y. Yang, F. Teng, P. Liu, M. Qi, C. Lv, J. Li, X. Zhang, and Z. He, "Haven: Hallucination-mitigated llm for verilog code generation aligned with hdl engineers," in *2025 Design, Automation & Test in Europe Conference (DATE)*. IEEE, 2025, pp. 1–7.
- [8] M. Akyash and H. M. Kamali, "Evolutionary large language models for hardware security: A comparative survey," in *Proceedings of the great lakes symposium on VLSI 2024*, 2024, pp. 496–501.
- [9] S. Feng, W. Shi, Y. Wang, W. Ding, V. Balachandran, and Y. Tsvetkov, "Don't hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration," *arXiv preprint arXiv:2402.00367*, 2024.
- [10] J. Song, X. Wang, J. Zhu, Y. Wu, X. Cheng, R. Zhong, and C. Niu, "Raghat: A hallucination-aware tuning pipeline for llm in retrieval-augmented generation," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, 2024, pp. 1548–1558.
- [11] S. Wu, H. Fei, L. Pan, W. Y. Wang, S. Yan, and T.-S. Chua, "Combating multimodal llm hallucination via bottom-up holistic reasoning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 8, 2025, pp. 8460–8468.
- [12] K. Liao, T. Wang, and F. Yu, "Qse: Mitigating llm hallucinations through query-adaptive saliency-localized activation editing," in *International Joint Conference on Artificial Intelligence*. Springer, 2025, pp. 56–73.
- [13] Y. Qiu, Z. Zhao, Y. Ziser, A. Korhonen, E. M. Ponti, and S. Cohen, "Spectral editing of activations for large language model alignment," *Advances in Neural Information Processing Systems*, vol. 37, pp. 56958–56987, 2024.
- [14] O. Skean, M. R. Arefin, Y. LeCun, and R. Shwartz-Ziv, "Does representation matter? exploring intermediate layers in large language models," *arXiv preprint arXiv:2412.09563*, 2024.
- [15] O. Skean, M. R. Arefin, D. Zhao, N. Patel, J. Naghiyev, Y. LeCun, and R. Shwartz-Ziv, "Layer by layer: Uncovering hidden representations in language models," *arXiv preprint arXiv:2502.02013*, 2025.
- [16] S. Zhang, T. Yu, and Y. Feng, "Truthx: Alleviating hallucinations by editing large language models in truthful space," *arXiv preprint arXiv:2402.17811*, 2024.
- [17] K. Tasnia, A. Garcia, T. Farheen, and S. Rahman, "Veriopt: Ppa-aware high-quality verilog generation via multi-role llms," *arXiv preprint arXiv:2507.14776*, 2025.
- [18] Z. Pei, H.-L. Zhen, M. Yuan, Y. Huang, and B. Yu, "Betternv: Controlled verilog generation with discriminative guidance," *arXiv preprint arXiv:2402.03375*, 2024.
- [19] Y. Zhao, D. Huang, C. Li, P. Jin, Z. Nan, T. Ma, L. Qi, Y. Pan, Z. Zhang, R. Zhang *et al.*, "Codev: Empowering llms for verilog generation through multi-level summarization," *arXiv preprint arXiv:2407.10424*, 2024.
- [20] M. Liu, Y.-D. Tsai, W. Zhou, and H. Ren, "Craftrtl: High-quality synthetic data generation for verilog code models with correct-by-construction non-textual representations and targeted code repair," *arXiv preprint arXiv:2409.12993*, 2024.
- [21] C.-T. Ho, H. Ren, and B. Khailany, "Verilogcoder: Autonomous verilog coding agents with graph-based planning and abstract syntax tree (ast)-based waveform tracing tool," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 1, 2025, pp. 300–307.
- [22] Y. Zhao, H. Zhang, H. Huang, Z. Yu, and J. Zhao, "Mage: A multi-agent engine for automated rtl code generation," *arXiv preprint arXiv:2412.07822*, 2024.
- [23] M. Akyash, K. Azar, and H. Kamali, "Rtl++: Graph-enhanced llm for rtl code generation," *arXiv preprint arXiv:2505.13479*, 2025.
- [24] C. Deng, Y.-D. Tsai, G.-T. Liu, Z. Yu, and H. Ren, "Scalertl: Scaling llms with reasoning data and test-time compute for accurate rtl code generation," *arXiv preprint arXiv:2506.05566*, 2025.
- [25] N. Wang, B. Yao, J. Zhou, Y. Hu, X. Wang, Z. Jiang, and N. Guan, "Large language model for verilog generation with code-structure-guided reinforcement learning," in *2025 IEEE International Conference on LLM-Aided Design (ICLAD)*. IEEE, 2025, pp. 164–170.
- [26] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, "Evaluating large language models trained on code," *arXiv preprint arXiv:2107.03374*, 2021.
- [27] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, and C. Xiong, "Codegen: An open large language model for code with multi-turn program synthesis," *arXiv preprint arXiv:2203.13474*, 2022.
- [28] X. Yao, Y. Wang, X. Li, Y. Lian, R. Chen, L. Chen, M. Yuan, H. Xu, and B. Yu, "Rtlrewriter: Methodologies for large models aided rtl code optimization," in *Proceedings of the 43rd IEEE/ACM International Conference on Computer-Aided Design*, 2024, pp. 1–7.
- [29] X. Yao, H. Li, T. H. Chan, W. Xiao, M. Yuan, Y. Huang, L. Chen, and B. Yu, "Hdldebugger: Streamlining hdl debugging with large language models," *ACM Transactions on Design Automation of Electronic Systems*, 2024.
- [30] N. Mashnoor, M. Akyash, H. Kamali, and K. Azar, "Llm-ift: Llm-powered information flow tracking for secure hardware," in *2025 IEEE 43rd VLSI Test Symposium (VTS)*, 2025, pp. 1–5.
- [31] H. Wu, Z. He, X. Zhang, X. Yao, S. Zheng, H. Zheng, and B. Yu, "Chateda: A large language model powered autonomous agent for eda," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2024.
- [32] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin *et al.*, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *ACM Transactions on Information Systems*, vol. 43, no. 2, pp. 1–55, 2025.
- [33] K. Li, O. Patel, F. Viégas, H. Pfister, and M. Wattenberg, "Inference-time intervention: Eliciting truthful answers from a language model," *Advances in Neural Information Processing Systems*, vol. 36, pp. 41451–41530, 2023.
- [34] C. Burns, H. Ye, D. Klein, and J. Steinhardt, "Discovering latent knowledge in language models without supervision," *arXiv preprint arXiv:2212.03827*, 2022.
- [35] Z. Chen, X. Sun, X. Jiao, F. Lian, Z. Kang, D. Wang, and C.-Z. Xu, "Truth forest: toward multi-scale truthfulness in large language models through intervention without tuning," in *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence*, ser. AAAI'24/IAAI'24/EAAI'24. AAAI Press, 2024. [Online]. Available: <https://doi.org/10.1609/aaai.v38i19.30087>
- [36] M. Liu, N. Pinckney, B. Khailany, and H. Ren, "Verilogeval: Evaluating large language models for verilog code generation," in *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. IEEE, 2023, pp. 1–8.