

PLUTO: A BENCHMARK FOR EVALUATING EFFICIENCY OF LLM-GENERATED HARDWARE CODE

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) are increasingly used to automate hardware design tasks, including the generation of Verilog code. While early benchmarks focus primarily on functional correctness, efficient hardware design demands additional optimization for synthesis metrics such as area, delay, and power. Existing benchmarks fall short in evaluating these aspects comprehensively: they often lack optimized baselines or testbenches for verification. To address these gaps, we present Pluto, a benchmark and evaluation framework designed to assess the efficiency of LLM-generated Verilog designs. Pluto presents a comprehensive evaluation set of 114 problems with self-checking testbenches and multiple Pareto-optimal reference implementations. Experimental results show that state-of-the-art LLMs can achieve high functional correctness, reaching 78.3% at pass@1, but their synthesis efficiency still lags behind expert-crafted implementations, with area efficiency of 63.8%, delay efficiency of 65.9%, and power efficiency of 64.0% at eff@1. This highlights the need for efficiency-aware evaluation frameworks such as Pluto to drive progress in hardware-focused LLM research.

1 INTRODUCTION

Large Language Models (LLMs) are beginning to reshape hardware design by automating key steps in hardware design workflows, including Verilog code generation Thakur et al. (2023a;b); Liu et al. (2023a), optimization Yao et al. (2024); Guo & Zhao (2025), verification Qiu et al. (2024a), debugging Tsai et al. (2024), high-level synthesis Xiong et al. (2024), and post-synthesis metric estimation Abdelatty et al. (2025). While these advances highlight the potential of LLMs in hardware design, most research has focused on functional correctness of generated designs, with little attention to design quality metrics such as area, delay, and power.

In hardware design, the quality of Verilog code is not determined solely by functional correctness. Designs typically undergo logic synthesis, where Verilog code is mapped to gate-level implementations in a target technology. This process exposes critical efficiency metrics, such as silicon area, timing delay, and power consumption, that directly impact manufacturability and performance. Unlike software code, where correctness and execution speed often suffice, hardware code quality is inherently tied to these post-synthesis metrics.

In order to evaluate the functional correctness of LLM-generated Verilog code, several benchmarks have been proposed including VerilogEval Liu et al. (2023a) and RTLLM Lu et al. (2024). Recent efforts, including RTLRewriter Yao et al. (2024), ResBench Guo & Zhao (2025), GenBen Wan et al. (2025), and TuRTL Garcia-Gasulla et al. (2025), have begun to evaluate quality of LLM-generated hardware code in terms of post-synthesis metrics. However, these benchmarks face key limitations:

- **Absence of Optimal Ground Truth Solutions** True efficiency should be measured against implementations that are explicitly optimized for specific objectives such as silicon area, delay, or power consumption. Prior studies rely on canonical solutions from VerilogEval and RTLLM as reference solutions. Our analysis shows that these solutions are not the most optimal in terms of post-synthesis metrics.
- **Lack of Clock Latency Agnostic Testbenches** Many common optimization patterns, such as register pipelining, resource sharing, or FSM restructuring, introduce variations in clock-cycle latency between the optimized and unoptimized designs. To support fair evaluation,

testbenches must be self-checking and tolerant of different latency requirements. Existing benchmarks, however, assume identical latency between the reference model and the design under test, making them unsuitable for efficiency benchmarking.

In order to address these limitations, we introduce *Pluto*, the first benchmark designed to evaluate both *functional correctness* and *synthesis efficiency* of LLM-generated Verilog code. Our contributions are as follows:

- **Per-Metric Ground Truth Optimal Solutions.** We provide a suite of 114 problems where each is optimized for area, delay, and power separately, yielding Pareto-front optimal solutions. Our ground truth solutions are significantly more efficient than canonical solutions in RTLLM and VerilogEval.
- **Optimization-Aware Testbenches.** Each problem is accompanied by clock-cycle agnostic testbenches that accommodate varying latency requirements, ensuring robust evaluation of different optimization patterns.
- **Comprehensive Evaluation.** We adapt the $eff@k$ metric introduced in Qiu et al. (2024b) to measure the efficiency of hardware designs. Our extended metric is a three-dimensional vector that evaluates LLM-generated code across multiple objectives: area, delay and power.

2 RELATED WORK

Software Code Benchmarks Large Language Models (LLMs) have been extensively studied for code generation across both software and hardware domains, with most early benchmarks focusing primarily on functional correctness rather than efficiency. In software, works such as Mercury Du et al. (2024) and ENAMEL Qiu et al. (2024b) move beyond correctness to explicitly evaluate run-time efficiency of LLM-generated programs. The Mercury Du et al. (2024) benchmark contains LeetCode style problems. Each problem is accompanied by an expert-written solution that represents the most optimal implementation in terms of run-time efficiency. ENAMEL Qiu et al. (2024b) also introduces a Python benchmark to evaluate the run-time efficiency of LLM-generated code.

Hardware Code Benchmarks In hardware design, early work on LLM-generated Verilog emphasized functional correctness. VerilogEval Liu et al. (2023a) only evaluates whether the LLM generated code passes the testbench check, while RTLLM Lu et al. (2024) additionally checks if the generated code is synthesizable. More recent efforts have shifted toward assessing and improving the efficiency of LLM-generated designs, which can be categorized into two main categories: *Specifications-to-Efficient-Verilog* where the LLM is tasked with translating natural language instruction to optimized Verilog code directly, and *Unoptimized-Verilog-to-Optimized-Verilog*, where the LLM is tasked with rewriting an unoptimized Verilog code to optimized Verilog code.

In the *Specifications-to-Efficient-Verilog* formulation, the LLM is prompted with a natural language problem description and directly generates optimized Verilog. Benchmarks, such as GenBen Wan et al. (2025), TuRTL Garcia-Gasulla et al. (2025), evaluate these generations in functional correctness, synthesizability, and post-synthesis metrics such as area, delay, and power. However, it relies on VerilogEval problems as ground truth. These reference designs are not necessarily optimized for power, performance, or area, and thus do not represent true Pareto-optimal solutions. ResBench Guo & Zhao (2025) also does not define any gold-standard or reference-optimal implementations, which makes it difficult to quantitatively assess how close the generated solutions are to ideal results.

The *Unoptimized-Verilog-to-Optimized-Verilog* setting provides the LLM with a functionally correct but unoptimized Verilog implementation and asks it to produce a more efficient version. RTL-Rewriter Yao et al. (2024) enhances this with retrieval-augmented generation and feedback through the synthesis loop. However, RTLRewriter lacks associated testbenches, making it unsuitable for assessing the functional correctness of the generated code.

As summarized in Table 1, *Pluto* is the first benchmark to offer per-metric optimization, providing separate expert-optimized reference designs for area, delay, and power. This enables targeted, metric-specific evaluation of LLMs, an aspect missing from prior benchmarks.

Table 1: Comparison of prior software and hardware code generation benchmarks. *Pluto* addresses key limitations by enabling metric-specific optimization with three reference implementations per problem, each optimized for area, delay, or power.

Benchmark	Language	Functionality	Synthesizability	Efficiency	Per-Metric Optimisation	Tasks
HumanEval	Python	✓	—	×	×	164
Mercury	Python	✓	—	✓	×	256
ENAMEL	Python	✓	—	✓	×	142
VerilogEval	Verilog	✓	×	×	×	156
RTLLe	Verilog	✓	✓	×	×	30
RTLRewriter	Verilog	×	✓	✓	×	95
ResBench	Verilog	✓	✓	×	×	56
GenBen	Verilog	✓	✓	×	×	300
TuRTLe	Verilog	✓	✓	×	×	223
CVDP	Verilog	✓	✓	×	×	783
Pluto (Ours)	Verilog	✓	✓	✓	✓	114

3 PLUTO BENCHMARK

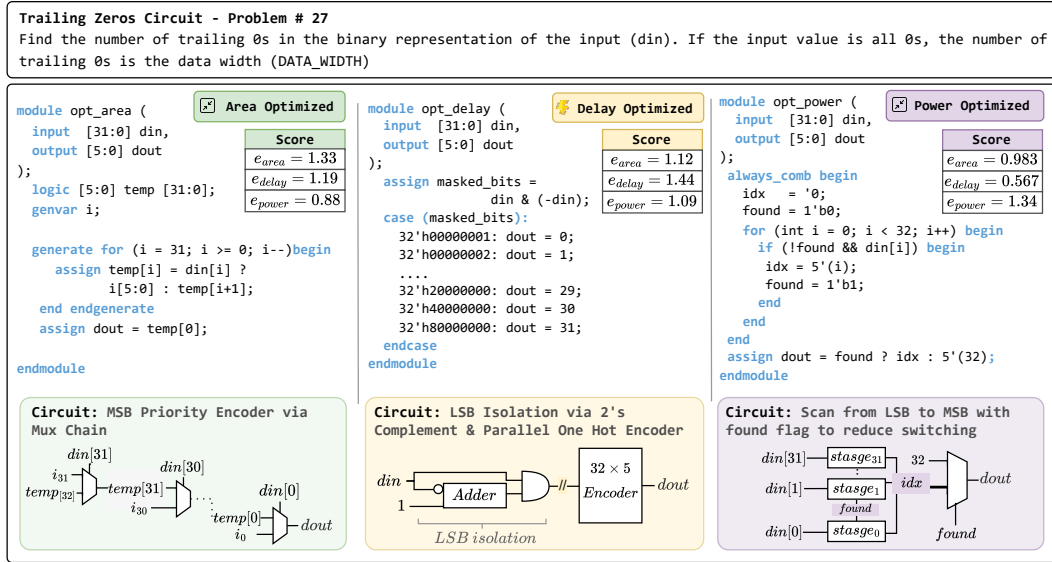


Figure 1: Overview of the *Pluto* benchmark on the trailing zeros detection task. We show three reference implementations optimized for different synthesis metrics compared to the unoptimized baseline: (left) area, using a mux-based priority encoder, reducing area by 33%; (center) delay, using an LSB isolation circuit with a parallel one-hot encoder, reducing delay by 44%; and (right) power, using an LSB-to-MSB scanning method with early termination, reducing total power by 34%. See Appendix. A.1 for unoptimized baseline and self-checking testbench.

3.1 DATA CONSTRUCTION

To enable a comprehensive evaluation of synthesis efficiency for LLM-generated hardware code, we construct the *Pluto* evaluation set, which contains a diverse collection of high-quality digital design problems spanning a broad range of difficulties. Specifically, we curated 114 problems from various publicly available sources, including open-source hardware projects, educational platforms such as ChipDev ChipDev (2025), a LeetCode-inspired platform for practicing Verilog coding, and prior benchmark suites such as RTLRewriter Yao et al. (2024), RTLLe Lu et al. (2024), and Ver-

ilogEval Liu et al. (2023b). Each problem is specified by a high-level description outlining the functional requirements, together with a baseline unoptimized Verilog implementation.

The problem set covers a wide spectrum of tasks in digital logic design, ranging from arithmetic units and control circuits to sequential state machines. To systematically capture variation in design complexity, we adopt ChipDev’s difficulty annotations and classify problems into three levels: easy, medium, and hard. These labels reflect the intrinsic challenge of translating the textual description into a correct Verilog implementation, thereby providing a principled way to distinguish between problems of different complexity.

Importantly, the resulting collection balances accessibility with challenge: many problems that appear straightforward can nonetheless expose substantial differences in synthesis efficiency depending on the optimization strategies applied. The diverse composition of easy, medium, and hard tasks therefore enables a nuanced assessment of an LLM’s ability to generate synthesis-efficient Verilog under varying constraints. In total, the 114 selected problems provide a representative and scalable testbed for benchmarking LLM-based Verilog efficiency.

Each problem instance in the *Pluto* benchmark includes the following components:

- **Prompt:** A natural language description of the hardware design task intended to guide the LLM-generation.
- **Module Header:** A fixed interface shared across all versions of the Verilog module to ensure consistency and comparability.
- **Unoptimized Verilog Code:** A baseline implementation used as the reference for testing.
- **Optimized Verilog Code:** Three distinct implementations with tradeoffs, each optimized by hand using design experts for a single metric: area, delay, or power.
- **Testbench:** A manually crafted, fully self-checking testbench that verifies functional equivalence between the unoptimized and any optimized design. These testbenches ensure full input space coverage and flag any mismatches during simulation. For sequential circuits, testbenches are clock-cycle agnostic, supporting latency differences introduced by optimizations such as pipelining or resource sharing.

All components in the evaluation set are manually developed. This ensures high quality and guarantees that the LLM under evaluation has not previously encountered any part of the dataset during training. In particular, the testbenches and optimized code serve as held-out ground truth references, providing an unbiased benchmark for assessing the efficiency and correctness of LLM-generated Verilog designs.

To illustrate the structure of problems in the *Pluto* evaluation set, Figure 1 presents the example of a trailing zeros detection circuit, categorized as an easy problem, along with its three metric-specific optimizations. As shown, each optimization achieves peak efficiency in its targeted metric, while performance in the remaining two metrics declines. This behavior emphasizes the inherent trade-offs across design objectives in hardware design and highlights the necessity of metric-specific optimization strategies.

3.2 OPTIMIZATION WORKFLOW

Each unoptimized design in the *Pluto* set is further refined through manual optimization by expert engineers to generate three distinct versions optimized separately for area, delay, and power. This workflow follows a systematic process that ensures both the correctness and the efficiency of the resulting designs. After applying metric-specific transformations, each optimized circuit is rigorously verified for functional correctness using Icarus Verilog Williams et al. (2002), supported by robust self-checking testbenches that guarantee equivalence with the unoptimized baseline. The optimized versions are then synthesized to confirm that improvements translate into measurable gains in area, timing, or power, thereby providing reliable performance baselines against which LLM-generated designs can be evaluated.

To understand how these efficiency gains are achieved, we visualize the optimization strategies applied across the dataset in appendix A.2. The strategies vary significantly depending on the target metric. For area, arithmetic optimizations and logic simplification are most commonly employed,

and FSM restructuring plays an important role in reducing redundant states and transitions. Delay improvements rely heavily on exploiting parallelism and restructuring control logic, often complemented by logic simplification and pipelining techniques that shorten the critical path. For power, it's reducing switching activity through register and logic optimizations, supported by techniques such as operand isolation, and clock gating to further suppress unnecessary toggling.

The distribution of strategies reveals that no single optimization technique dominates across all objectives. Instead, engineers select strategies tailored to the specific metric, reflecting the trade-offs inherent in digital design. As shown in Figure 2, this process results in consistent improvements across **all 114** designs, with average reductions of 19.19% ($SD=18.99\%$) in area, 21.96% ($SD=20.99\%$) in delay, and 22.55% ($SD=21.65\%$) in power. This highlights the importance of metric-specific approaches and provide a robust baseline for evaluating LLM-generated hardware code efficiency.

To further illustrate the impact of expert-driven optimization, Table 2 presents representative examples drawn from both RTLLM and VerilogEval. These case studies highlight how different strategies, such as arithmetic unit sharing, FSM encoding choices, and counter-based control logic, translate into concrete improvements across area, delay, and power. As shown, across both VerilogEval and RTLLM, expert-optimized designs consistently outperform baseline implementations. In particular, for RTLLM problems our expert-written solutions achieve average improvements of 18.75% ($SD=14.55\%$) in area, 22.75% ($SD=20.99\%$) in delay, and 20.43% ($SD=21.65\%$) in power compared to their canonical solutions. For VerilogEval problems, the improvements average 10.46% ($SD=14.40\%$) in area, 10.33% ($SD=15.10\%$) in delay, and 13.61% ($SD=18.86\%$) in power.

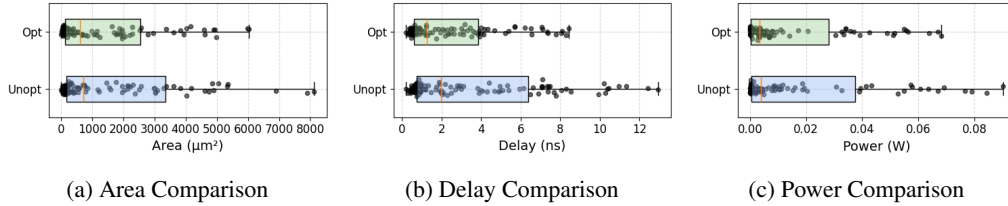


Figure 2: Distribution of area, delay, and power across **all 114 tasks** in *Pluto* before and after manual metric-specific optimizations.

Table 2: A sample of benchmark problems from *Pluto* dataset. Our expert-optimized solutions (area, delay, power) are significantly more efficient than the baseline benchmark implementations. See Appendix A.3 for full problem implementation.

ID	Source	Problem Description	Benchmark Solution	Expert Solution (Ours)
#60	RTLLM	ALU for a 32-bit MIPS-ISA CPU with operations like {ADD, SUB, AND, OR, XOR, SLT, shifts, LUI}.	ALU implementation with case statement and parameterized op-codes.	Area ↓ 26%: Shared adder for arithmetic, simplified flag logic, and operand reuse Delay ↓ 26%: Parallel datapaths with one-hot muxing Power ↓ 4%: Operand gating and early zeroing for large shifts to cut switching activity
#68	RTLLM	FSM detecting the sequence 10011 on a serial input stream with support for continuous and overlapping detection	States binary encoded, sequential next-state logic, and registered Mealy output	Area ↓ 32%: Casez-based transitions and direct output Delay ↓ 17%: One-hot state encoding with pre-decoded inputs and Moore-style output Power ↓ 23%: Compact binary encoding and casez-based transitions to cut toggling
#87	VerilogEval	Module controls a shift register enable signal, <i>shift_ena</i> asserted for 4 clock cycles on reset, then remain 0 until the next reset	Uses explicit states with next-state logic to drive <i>shift_ena</i>	Area ↓ 47%: Minimized register width (2-bit counter) and compact comparator logic Delay ↓ 37%: Wider counter (3 bits) to simplify comparison and reduce logic depth on the critical path Power ↓ 46%: Small counter reused
#104	VerilogEval	Conway's Game of Life with a 16×16 toroidal grid: each cell updates based on neighbor counts n (live if $n = 3$ or $n = 2 \ \& \ \text{alive}$)	Straightforward RTL with per-cell neighbor recomputation and sequential summing	Area ↓ 19%: Shared neighbor computations across rows, bitwise rotations for wraparound, less duplicate summations Delay ↓ 37%: Parallel neighbor summation with carry-save adder tree and direct decode for 2 and 3 Power ↓ 36%: Reduced toggling via computation reuse

3.3 EFFICIENCY METRICS

We use the $pass@k$ Liu et al. (2023a) for measuring the functional correctness of LLM-generated Verilog code. The $pass@k$ metric, defined in appendix. A.4, measures the percentage of problems for which at least one of the top- k generated samples passes the self-checking testbench.

To evaluate the synthesis efficiency of functionally correct samples, we adapt the $eff@k$ introduced in Qiu et al. (2024b) to Verilog code. First, we introduce the efficiency score $e_{i,j}$, defined in Eq. 1, which quantifies how close an LLM-generated design is to optimal ground truth implementation. In this equation, $\hat{R}_{i,j}$ denotes the reported synthesis metric (e.g., area, delay, or power) for the j -th sample of problem i , $T_{i,j}$ denotes an upper bound beyond which the design is considered inefficient, and $R_{i,j}$ denotes the optimal (lowest) known reference value for that metric. A score of 1 indicates that the sample exactly matches the optimal reference, while a score of 0 indicates that it exceeds the acceptable threshold or is functionally incorrect.

$$e_{i,j} = \begin{cases} \frac{\max(0, T_{i,j} - \hat{R}_{i,j})}{T_{i,j} - R_{i,j}}, & \text{if } n_{i,j} \text{ is correct} \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

$$eff@k = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{J \subseteq \{1, \dots, n\}, |J|=k} \left[\max_{j \in J} e_{i,j} \right] \\ = \frac{1}{N} \sum_{i=1}^N \sum_{r=k}^n \frac{\binom{r-1}{k-1}}{\binom{n}{k}} e_{i,(r)}. \quad (2)$$

We then use the efficiency score $e_{i,j}$ for computing the $eff@k$, defined in Eq. 2 as the average of the best (i.e., highest) efficiency scores among the top- k functionally correct samples for each problem. We use the unbiased estimator introduced in Qiu et al. (2024b) for computing $eff@k$ which computes the expectation value over a random subset J of code samples with size K .

4 EVALUATION RESULTS

We evaluate 18 large language models (LLMs) using our *Pluto* benchmark, which includes proprietary LLMs, general-purpose foundation models, code-specialized models, and Verilog-tuned models. To comprehensively assess efficiency-aware generation, we consider the two problem formulations in *Pluto*: translating unoptimized Verilog code into optimized implementations, and generating optimized code directly from natural-language specifications. For the first problem formulation, only instruction-tuned models are evaluated, as code completion models generally reproduce the unoptimized code without meaningful improvements.

4.1 MAIN RESULTS

Table 3 (a) reports the $pass@k$ and $eff@k$ metrics for the first problem formulation, where the task is to re-write unoptimized Verilog into more efficient implementations. Several trends emerge. First, in terms of functional correctness ($pass@k$), domain-tuned models such as VeriThoughts-Inst-7B and RTLCoder-DeepSeek-V1 achieve performance comparable to much larger foundational models like DeepSeek-Chat, demonstrating the benefit of Verilog-specific training. However, in terms of synthesis efficiency ($eff@k$), all models exhibit a noticeable drop relative to their $pass@k$ scores. This gap underscores a common limitation: while LLMs can generate functionally correct Verilog, they struggle to match the Pareto-efficient expert baselines across area, delay, and power.

Table 3 (b) reports the $pass@k$ and $eff@k$ metrics for the second problem formulation, where models are tasked with translating natural language specifications into optimized Verilog implementations. This task is more challenging, and as a result, both $pass@k$ and $eff@k$ scores are consistently lower across all models. Similar to the first formulation, all models also exhibit lower $eff@k$ values compared to their corresponding $pass@k$ scores, underscoring the persistent difficulty of generating designs that are not only functionally correct but also synthesis-efficient. However, the relative gap between $pass@k$ and $eff@k$ is smaller in this setting compared to the first formulation. This is because specification-to-RTL translation is substantially harder: models often struggle to produce functionally correct code in the first place, which suppresses both correctness and efficiency scores.

Table 3: Evaluation results using *Pluto* for two problem formulations: **P1: Unoptimized-Verilog-to-Optimized-Verilog** and **P2: Specifications-to-Optimized-Verilog**. $pass@k$ measures functional correctness, while $eff@k$ measures efficiency across area, delay, and power.

(a) **P1: Unoptimized-Verilog-to-Optimized-Verilog**

Model	pass@1	pass@5	pass@10	eff@1			eff@5			eff@10		
				Area	Delay	Power	Area	Delay	Power	Area	Delay	Power
GPT-3.5	0.325	0.517	0.594	0.271	0.296	0.282	0.462	0.491	0.450	0.540	0.568	0.520
GPT-4o-mini	0.506	0.705	0.751	0.469	0.476	0.467	0.662	0.677	0.639	0.714	0.744	0.687
DeepSeek-Chat	0.612	0.802	0.860	0.586	0.599	0.601	0.776	0.795	0.794	0.839	0.862	0.846
Llama-3.3-70B-Instruct	0.473	0.701	0.757	0.446	0.462	0.429	0.662	0.696	0.662	0.707	0.760	0.735
Llama-3.1-8B-Instruct	0.160	0.432	0.567	0.127	0.156	0.145	0.358	0.437	0.384	0.494	0.584	0.505
Mistral-7B-Instruct-v0.2	0.106	0.318	0.453	0.078	0.094	0.100	0.244	0.296	0.301	0.358	0.446	0.427
Mistral-8x7B-v0.1	0.255	0.520	0.652	0.217	0.231	0.210	0.462	0.487	0.447	0.593	0.630	0.561
starcoder2-15b-instruct-v0.1	0.659	0.960	0.988	0.611	0.633	0.591	0.879	0.924	0.871	0.913	0.952	0.904
CodeLlama-70b-Instruct-hf	0.576	0.905	0.956	0.522	0.541	0.523	0.842	0.876	0.824	0.903	0.925	0.878
DeepSeek-Coder-33B	0.783	0.963	0.997	0.638	0.659	0.640	0.902	0.927	0.883	0.942	0.960	0.927
Qwen2.5-Coder-7B-Inst	0.479	0.785	0.866	0.438	0.452	0.419	0.710	0.759	0.741	0.785	0.848	0.833
yang-z/CodeV-QC-7B	0.231	0.416	0.506	0.211	0.208	0.187	0.381	0.390	0.361	0.455	0.491	0.442
RTLCoder-DeepSeek-V1	0.532	0.854	0.915	0.471	0.495	0.468	0.774	0.789	0.757	0.843	0.850	0.809
VeriThoughts-Inst.-7B	0.611	0.797	0.854	0.540	0.560	0.524	0.708	0.740	0.702	0.763	0.785	0.765

(b) **P2: Specifications-to-Optimized-Verilog**

Model	pass@1	pass@5	pass@10	eff@1			eff@5			eff@10		
				Area	Delay	Power	Area	Delay	Power	Area	Delay	Power
GPT-3.5	0.239	0.395	0.471	0.225	0.235	0.225	0.373	0.390	0.381	0.439	0.469	0.468
GPT-4o-mini	0.391	0.533	0.591	0.360	0.377	0.363	0.475	0.520	0.495	0.532	0.572	0.551
DeepSeek-Chat	0.557	0.688	0.719	0.545	0.552	0.528	0.689	0.680	0.651	0.726	0.710	0.684
Llama-3.3-70B-Instruct	0.363	0.541	0.594	0.348	0.345	0.342	0.515	0.533	0.510	0.564	0.602	0.557
Llama-3.1-8B-Instruct	0.087	0.224	0.301	0.075	0.073	0.081	0.189	0.184	0.204	0.270	0.251	0.279
Mistral-7B-Instruct-v0.2	0.030	0.108	0.164	0.024	0.015	0.024	0.078	0.067	0.088	0.112	0.117	0.134
Mistral-8x7B-v0.1	0.082	0.176	0.222	0.082	0.068	0.079	0.172	0.131	0.163	0.202	0.166	0.211
starcoder2-15b-instruct-v0.1	0.243	0.454	0.512	0.226	0.249	0.220	0.429	0.466	0.409	0.489	0.525	0.458
CodeLlama-70b-Instruct-hf	0.212	0.446	0.532	0.202	0.207	0.194	0.418	0.440	0.437	0.498	0.535	0.524
DeepSeek-Coder-33B	0.257	0.429	0.482	0.231	0.254	0.246	0.387	0.424	0.421	0.446	0.490	0.468
Qwen2.5-Coder-7B-Inst	0.164	0.324	0.389	0.158	0.162	0.148	0.307	0.319	0.295	0.366	0.375	0.357
code-gen-verilog-16b	0.068	0.200	0.289	0.069	0.064	0.058	0.188	0.175	0.188	0.268	0.253	0.272
yang-z/CodeV-CL-7B	0.265	0.485	0.553	0.233	0.243	0.237	0.432	0.455	0.436	0.493	0.536	0.511
yang-z/CodeV-QC-7B	0.260	0.458	0.529	0.223	0.229	0.222	0.420	0.415	0.410	0.497	0.480	0.478
yang-z/CodeV-All-QC	0.175	0.317	0.374	0.150	0.162	0.144	0.297	0.304	0.256	0.368	0.379	0.284
RTLCoder-DeepSeek-V1	0.203	0.400	0.480	0.177	0.199	0.184	0.345	0.404	0.361	0.417	0.499	0.430
RTLCoder-Mistral	0.199	0.347	0.418	0.185	0.188	0.188	0.316	0.332	0.329	0.381	0.411	0.395
VeriThoughts-Inst.-7B	0.216	0.336	0.398	0.211	0.206	0.207	0.316	0.330	0.329	0.367	0.394	0.393

4.2 ABLATION STUDIES

In addition to the Verilog code writing style, post-synthesis metrics are also influenced by external factors such as the synthesis tool employed, the target technology library, and the optimization sequence executed by the tool. To understand the robustness of our proposed benchmark and isolate the impact of these factors, we present two ablation studies that evaluate efficiency trends in the *Pluto* benchmark across different synthesis tools, optimization strategies, and technology libraries.

4.2.1 SYNTHESIS TOOL AND TECHNOLOGY AGNOSTICISM

In this experiment, we repeated synthesis runs for the three optimized reference implementations in *Pluto*, as well as the unoptimized baseline, using two distinct synthesis tools: Yosys Wolf et al. (2013), an open-source framework, and Cadence Genus cad, a commercial synthesis tool. To further evaluate generalizability, we also targeted two technology libraries representing different fabrication nodes: the SkyWater 130nm library Google and a 65nm TSMC library tsm. We then computed the efficiency score for each tool and library configuration by comparing each optimized implementation against the corresponding unoptimized baseline across area, delay, and power metrics. As shown in Figure 3, efficiency scores remain consistent across all synthesis tool and technology combinations. This demonstrates that *Pluto*’s optimization patterns deliver consistent tradeoffs across different synthesis tools and technology libraries.

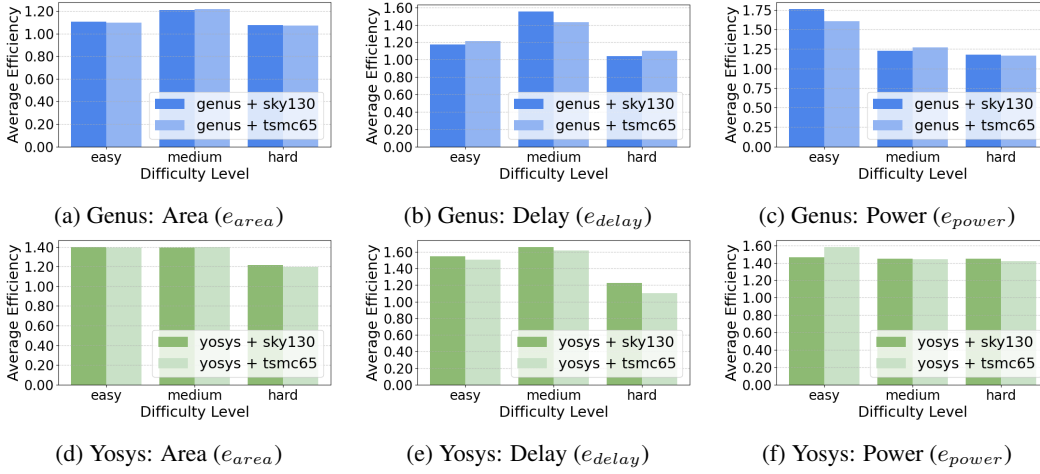


Figure 3: Efficiency scores for area, delay, and power across all benchmark problems, using both Cadence Genus and Yosys with different technology libraries. Results show consistent efficiency trends across synthesis tools and technologies.

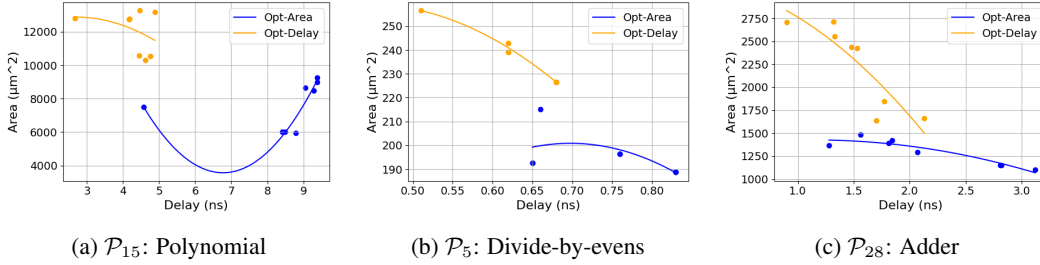


Figure 4: Area–delay tradeoffs of three problems in the *Pluto* benchmark under different synthesis strategies. Each strategy corresponds to a distinct sequence of ABC logic synthesis commands.

4.2.2 SYNTHESIS OPTIMIZATION STRATEGIES

We also examine how different synthesis optimization strategies influence the post-synthesis metrics of *Pluto*’s optimized implementations. Synthesis tools allow designers to specify optimization directives that steer the tool’s internal heuristics toward minimizing a particular metric while potentially sacrificing others. To study this effect, we synthesized selected problems from the *Pluto* benchmark under both area-optimized and delay-optimized optimization strategies. We used Yosys as our synthesis tool and targeted the SkyWater 130nm library. Within Yosys, logic optimization is carried out using the ABC framework Synthesis & Group (2024), which provides a collection of optimization heuristics that can be configured to emphasize different objectives such as area or delay minimization. Figure 4 illustrates the resulting Pareto fronts of area–delay trade-offs across three representative problems. As expected, delay-optimized code consistently achieves superior timing performance at the expense of larger area, whereas area-optimized code achieves lower area but incurs higher delays. These results confirm that synthesis settings primarily shift designs along the area–delay curve, while coding style remains the dominant factor, validating *Pluto*’s ability to capture design efficiency independent of synthesis optimization settings.

5 MULTI-OBJECTIVE OPTIMIZATION

Pluto also supports benchmarking LLMs for multi-objective optimization across area, delay, and power. For each problem, the individually optimized design variants represent distinct Pareto-optimal points. Together, they form a Pareto front that captures the trade-offs between these de-

Table 4: Multi-objective performance comparison across metric combinations using the P2 problem formulation (Specifications-to-Optimized-Verilog). Best metric set per model is bolded.

Model	Area–Delay		Area–Power		Delay–Power		Area–Delay–Power	
	pass@10	eff@10	pass@10	eff@10	pass@10	eff@10	pass@10	eff@10
GPT-4o-mini	0.570	0.501	0.597	0.548	0.570	0.515	0.579	0.506
Llama3.3-70B-Instruct	0.649	0.615	0.605	0.577	0.640	0.600	0.605	0.567
Qwen2.5-Coder-7B-Instruct	0.395	0.358	0.395	0.364	0.421	0.393	0.430	0.402
Mixtral-8x7B	0.219	0.193	0.219	0.204	0.237	0.207	0.228	0.207
StarCoder2	0.500	0.465	0.526	0.499	0.526	0.502	0.465	0.440
CodeLlama-70B	0.526	0.490	0.553	0.529	0.535	0.499	0.535	0.512
DeepSeek-Coder-33B	0.474	0.444	0.518	0.480	0.491	0.463	0.465	0.439
yang-z/CodeV-QC-7B	0.526	0.474	0.526	0.472	0.544	0.488	0.509	0.472
RTLCoder-DeepSeek-V1	0.483	0.409	0.465	0.437	0.439	0.399	0.474	0.450
VeriThoughts-Instruct-7B	0.386	0.368	0.404	0.389	0.421	0.411	0.342	0.323

sign metrics, enabling the evaluation of LLM-generated code for efficiency across multiple objectives. To generalize the efficiency score to multiple metrics, we consider a task-specific metric set $\mathcal{M} \subseteq \{\text{area, delay, power}\}$, where the size of \mathcal{M} determines the dimensionality of the objective (e.g., two metrics for area–delay optimization or all three for full PPA evaluation). For each metric $m \in \mathcal{M}$, let $e_{i,j}^{(m)}$ denote the single-metric efficiency score (defined in Eq. 1). We then compute a multi-objective efficiency score as a weighted combination defined in Eq. 3. Finally, we obtain the multi-objective eff@k by substituting $e_{i,j}^{\text{multi}}$ for $e_{i,j}$ in Eq. 2.

$$e_{i,j}^{\text{multi}} = \sum_{m \in \mathcal{M}} w_m e_{i,j}^{(m)}, \quad \sum_{m \in \mathcal{M}} w_m = 1, \quad w_m \geq 0. \quad (3)$$

Table 4 shows the performance of different LLMs on multi-objective optimization across varying metric combinations. For each metric set \mathcal{M} , we report the multi-objective eff@k scores with equal weighting ($w_m=1/|\mathcal{M}|$) across all metrics in the set. The results show that, unlike single-metric evaluation, there is no universally easiest metric combination. For example, GPT-4o-mini peaks on Area–Power, while the Llama3.3-70B-Instruct model achieves its highest score on Area–Delay. Several mid-sized models, such as Mixtral and StarCoder2, achieve their best performance on Delay–Power. Across most models, the lowest scores appear on the Area–Delay or Area–Delay–Power combinations, reflecting the inherently competing nature of these metrics. Overall, these results indicate that multi-objective PPA optimization remains a challenging problem, and different LLMs exhibit different strengths depending on the metric set.

6 FAILURE ANALYSIS AND INSIGHTS

While LLMs reliably produce functionally correct Verilog, their ability to optimize is uneven across metrics. Area optimization is comparatively tractable, since it often reduces to logic simplification or FSM re-encoding. These area optimizations represent common, syntactic, pattern-like edits that appear frequently in code corpora and programmer-annotated examples, making them more learnable for LLMs. By contrast, delay requires identifying and shortening the critical path, and power depends on subtle factors like switching activity and memory usage. Crucially, area transforms are often local (e.g., simplifying a logic expression), while power and delay optimization typically requires global reasoning across the entire design (e.g., pipeline balancing, critical path restructuring). Current LLMs, especially smaller ones, struggle with these larger-scope transformations. This difficulty is reflected in our quadrant analysis (Figure 5a and 8), where many delay- and power-optimized designs remain correct but fail to improve efficiency, whereas area optimizations succeed more often. Moreover, analysis of optimization strategies (Figure 5b) shows that the hardest transformations are register optimizations for delay, followed by resource sharing for power, and sequential restructuring for delay. In contrast, strategies tied to area are easier, aligning with our observation that area is the most accessible metric for LLMs.

Model scale and specialization strongly influence outcomes. Larger models (15B, 33B, 70B) capture richer patterns and propose alternative architectures, showing stronger generalization across

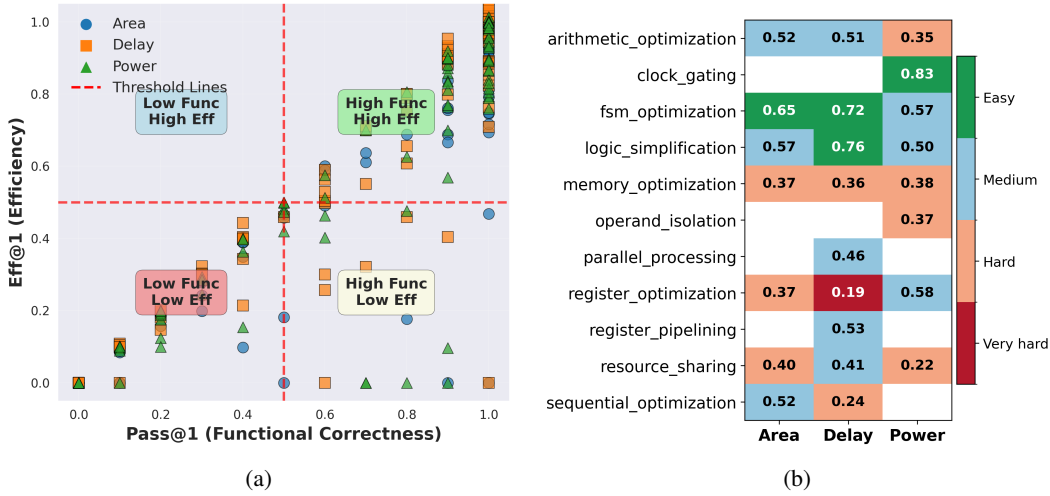


Figure 5: Failure mode analysis of optimization outcomes. (a) Quadrant plot showing the correlation between functional correctness ($Pass@1$) and synthesis efficiency ($Eff@1$) across area, delay, and power objectives. (b) Heatmap of optimization strategy difficulty across different optimization objectives area, delay, and power.

most optimization tactics. For example, larger code models perform more consistently on FSM optimization compared to smaller models, though they still struggle with register optimization and resource sharing. Models with explicit reasoning traces (e.g., DeepSeek, VeriThoughts) better decompose transformations and achieve stronger optimizations. Interestingly, register optimization lags behind other strategies across all DeepSeek variants (DeepSeek-Chat, DeepSeek-Coder, RTLCode-DeepSeek-V1), suggesting this tactic requires capabilities beyond what current architectures capture, regardless of scale or specialization. Domain-tuned models outperform code models, which in turn outperform general-purpose LLMs, showing the value of Verilog-specific pretraining. Notably, generalist chat models show weaker performance on power-optimization techniques like clock gating and operand isolation, which are less prevalent in general training datasets.

Finally, a fundamental limitation is that Verilog training data lacks efficiency labels. LLMs therefore default to surface-level pattern matching rather than structural reasoning, and without feedback or synthesis-in-the-loop, they cannot tell whether changes reduce gate count or lengthen the critical path. Completion-style models exacerbate this issue, often rephrasing the baseline instead of innovating, whereas instruction-tuned models attempt more substantive edits. These findings suggest that true progress will require efficiency-focused benchmarks such as Pluto to guide future advances.

7 CONCLUSION

In this paper, we introduced *Pluto*, a comprehensive benchmark designed to evaluate the synthesis efficiency of LLM-generated Verilog code. *Pluto* provides an evaluation set of 114 hardware design problems, each accompanied by three reference optimized implementations (targeting area, delay, and power), an unoptimized baseline, a self-checking testbench, and a natural language description. Experimental results show that while LLMs can achieve high functional correctness, reaching up to 78.3% at pass@1, their synthesis efficiency remains limited: area efficiency of 63.8%, delay efficiency of 65.9%, and power efficiency of 64.0% at eff@1 compared to expert-crafted designs. These findings highlight the importance of efficiency-aware benchmarks beyond correctness alone and highlights the current limitations of LLMs in hardware optimization.

REFERENCES

- Cadence Genus Synthesis Solution. https://www.cadence.com/en_US/home/tools/digital-design-and-signoff/synthesis/genus-synthesis.html.
- TSMC 65nm Technology Library (Proprietary). Accessed under license from TSMC.
- Manar Abdelatty, Jingxiao Ma, and Sherief Reda. Metrex: A benchmark for verilog code metric reasoning using llms. In *Proceedings of the 30th Asia and South Pacific Design Automation Conference*, pp. 995–1001, 2025.
- ChipDev. ChipDev: Hardware Interview Prep & Verilog Practice. <https://chipdev.io>, 2025. Accessed: 2025-01-01.
- Mingzhe Du, Anh Tuan Luu, Bin Ji, Qian Liu, and See-Kiong Ng. Mercury: A code efficiency benchmark for code large language models. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- Dario Garcia-Gasulla, Gokcen Kestor, Emanuele Parisi, Miquel Albert’i-Binimelis, Cristian Gutierrez, Razine Moundir Ghorab, Orlando Montenegro, Bernat Homs, and Miquel Moreto. Turtle: A unified evaluation of llms for rtl generation. *arXiv preprint arXiv:2504.01986*, 2025.
- Google. Skywater-pdk. <https://github.com/google/skywater-pdk>. Accessed: 2025-02-09.
- Ce Guo and Tong Zhao. Resbench: Benchmarking llm-generated fpga designs with resource awareness. *arXiv preprint arXiv:2503.08823*, 2025.
- Mingjie Liu, Nathaniel Pinckney, Brucek Khailany, and Haoxing Ren. Invited paper: Verilogeval: Evaluating large language models for verilog code generation. In *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*, pp. 1–8, 2023a. doi: 10.1109/ICCAD57390.2023.10323812.
- Mingjie Liu, Nathaniel Pinckney, Brucek Khailany, and Haoxing Ren. Verilogeval: Evaluating large language models for verilog code generation. In *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*, pp. 1–8. IEEE, 2023b.
- Yao Lu, Shang Liu, Qijun Zhang, and Zhiyao Xie. Rtlm: An open-source benchmark for design rtl generation with large language model. In *2024 29th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 722–727. IEEE, 2024.
- Ruidi Qiu, Grace Li Zhang, Rolf Drechsler, Ulf Schlichtmann, and Bing Li. Autobench: Automatic testbench generation and evaluation using llms for hdl design. In *Proceedings of the 2024 ACM/IEEE International Symposium on Machine Learning for CAD, MLCAD ’24*, New York, NY, USA, 2024a. Association for Computing Machinery. ISBN 9798400706998. doi: 10.1145/3670474.3685956. URL <https://doi.org/10.1145/3670474.3685956>.
- Ruizhong Qiu, Weiliang Will Zeng, James Ezick, Christopher Lott, and Hanghang Tong. How efficient is llm-generated code? a rigorous & high-standard benchmark. *arXiv preprint arXiv:2406.06647*, 2024b.
- Berkeley Logic Synthesis and Verification Group. Abc: A system for sequential synthesis and verification. Technical report, University of California, Berkeley, 2024. URL <https://people.eecs.berkeley.edu/~alanmi/abc/>.
- Shailja Thakur, Baleegh Ahmad, Zhenxing Fan, Hammond Pearce, Benjamin Tan, Ramesh Karri, Brendan Dolan-Gavitt, and Siddharth Garg. Benchmarking large language models for automated verilog rtl code generation. In *2023 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1–6. IEEE, 2023a.
- Shailja Thakur, Baleegh Ahmad, Hammond Pearce, Benjamin Tan, Brendan Dolan-Gavitt, Ramesh Karri, and Siddharth Garg. Verigen: A large language model for verilog code generation. *ACM Transactions on Design Automation of Electronic Systems*, 2023b.

YunDa Tsai, Mingjie Liu, and Haoxing Ren. Rtlfixer: Automatically fixing rtl syntax errors with large language models. In *IEEE/ACM Design Automation Conference (DAC'24)*, pp. 1–8, 2024. doi: 10.1109/ICCAD57390.2023.10323812.

Gwok-Waa Wan, Yubo Wang, SamZaak Wong, Jingyi Zhang, Mengnv Xing, Zhe Jiang, Nan Guan, Ying Wang, Ning Xu, Qiang Xu, and Xi Wang. Genben: A generative benchmark for LLM-aided design. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=gtVo4xcpFI>. Under review.

Stephen Williams et al. Icarus verilog: open-source verilog more than a year later. *Linux Journal*, 2002.

Clifford Wolf, Johann Glaser, and Johannes Kepler. Yosys—a free verilog synthesis suite. In *21st Austrian Workshop on Microelectronics (Austrochip)*, volume 97, 2013.

Chenwei Xiong, Cheng Liu, Huawei Li, and Xiaowei Li. Hlsplit: Llm-based high-level synthesis. *arXiv preprint arXiv:2408.06810*, 2024.

Xufeng Yao, Yiwen Wang, Xing Li, Yingzhao Lian, Ran Chen, Lei Chen, Mingxuan Yuan, Hong Xu, and Bei Yu. Rtlrewriter: Methodologies for large models aided rtl code optimization. *arXiv preprint arXiv:2409.11414*, 2024.

A APPENDIX

A.1 UNOPTIMIZED CODE AND TESTBENCH FOR PROBLEM #17 (TRAILING ZEROS) IN FIGURE 1

A.1.1 UNOPTIMIZED CODE

```

1 module unopt_model #(parameter
2   DATA_WIDTH = 32
3 ) (
4   input  [DATA_WIDTH-1:0] din,
5   output logic [$clog2(DATA_WIDTH):0] dout
6 );
7
8   logic [DATA_WIDTH-1:0] din_adj;
9   logic [$clog2(DATA_WIDTH):0] idx;
10
11   always_comb begin
12       idx = 0;
13       din_adj = din & (~din+1);
14       for (int i=0; i<DATA_WIDTH; i++) begin
15           idx += (din_adj[i]) ? i : 0;
16       end
17   end
18
19   assign dout = (din_adj == 0 ? DATA_WIDTH : din_adj == 1 ? 0 : idx);
20
21 endmodule

```

A.1.2 SELF-CHECKING TESTBENCH

```

1 `timescale 1 ps/1 ps
2
3 module tb();
4
5   reg clk = 0;
6   initial forever #5 clk = ~clk;
7
8   wire [5:0] dout_opt, dout_unopt;
9   reg [31:0] din;
10
11   integer errors = 0;
12   integer errortime = 0;
13   integer clocks = 0;
14   integer total_cycles = 200;
15
16   initial begin

```

```

648 17      $dumpfile("wave.vcd");
649 18      $dumpvars(1, clk, din, dout_opt, dout_unopt);
650 19
651 20      // Initialize din to avoid X values
652 21      din = 0;
653 22
654 23      // Generate random values for din
655 24      repeat(total_cycles) @(posedge clk) din = $random;
656 25  end
657 26
658 27  wire tb_match;
659 28  assign tb_match = (dout_opt === dout_unopt);
660 29
661 30  opt_model opt_model (
662 31      .din(din),
663 32      .dout(dout_opt)
664 33  );
665 34
666 35  unopt_model unopt_model (
667 36      .din(din),
668 37      .dout(dout_unopt)
669 38  );
670 39
671 40  always @(posedge clk) begin
672 41      clocks = clocks + 1;
673 42      if (!tb_match) begin
674 43          if (errors == 0) errortime = $time;
675 44          errors = errors + 1;
676 45      end
677 46
678 47      // Print the signals for debugging
679 48      $display("Time=%0t, Cycle=%0d, din=%h, dout_opt=%h, dout_unopt=%h, match=%b",
680 49          $time, clocks, din, dout_opt, dout_unopt, tb_match);
681 50
682 51      if (clocks >= total_cycles) begin
683 52          $display("Simulation completed.");
684 53          $display("Total mismatches: %d out of %d samples", errors, clocks);
685 54          $display("Simulation finished at %0d ps", $time);
686 55          $finish;
687 56      end
688 57  end
689 58
690 59  initial begin
691 60      #1000000
692 61      $display("TIMEOUT");
693 62      $finish();
694 63  end
695 64
696 65  endmodule

```

A.2 OPTIMIZATION STRATEGIES

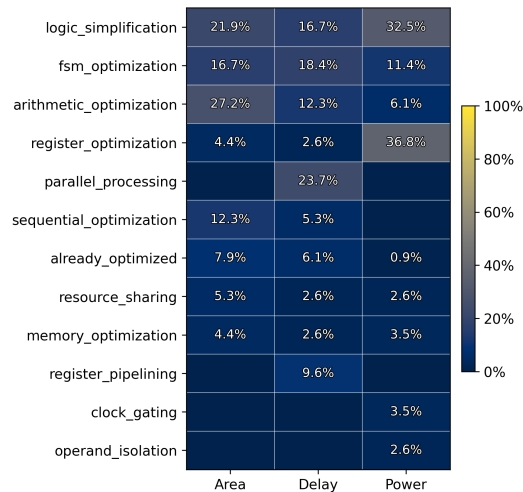


Figure 6: Optimization strategies employed for area, delay, and power improvements.

A.3 CODE OF EXAMPLE PROBLEMS IN TABLE 2

In the example problems shown in the appendix, the area- and power-optimized solutions coincided, as area-oriented designs also achieved the best power results, and vice versa. This overlap arises because common power-saving techniques, such as clock gating and operand isolation, were not applicable as some designs lacked a clock signal, while others did not include an enable signal. Consequently, explicit power-specific transformations could not be meaningfully applied. Moreover, in certain cases, power optimizations indirectly reduced area, further reinforcing the convergence of the two objectives into a single optimized implementation.

A.3.1 PROBLEM #60: RTLLM ALU

Problem description: Implement a 32-bit Arithmetic Logic Unit (ALU) for a MIPS-ISA CPU. The ALU takes two 32-bit operands (a and b) and a 6-bit control signal (aluc) that specifies which operation to perform. Based on this control signal, the ALU produces a 32-bit result (r) and several status outputs: zero indicates whether the result is zero, carry flags if a carry occurred, negative shows if the result is negative, overflow signals arithmetic overflow, and flag is used for set-less-than instructions (slt and sltu). The module supports arithmetic, logical, shift, and immediate load operations defined by specific opcodes (e.g., ADD, SUB, AND, OR, XOR, SLT, LUI).

Benchmark solution: ALU implementation with case statement and parameterized opcodes.

```

721 module unopt_model(
722     input [31:0] a,
723     input [31:0] b,
724     input [5:0] aluc,
725     output [31:0] r,
726     output zero,
727     output carry,
728     output negative,
729     output overflow,
730     output flag
731 );
732
733 parameter ADD = 6'b100000;
734 parameter ADDU = 6'b100001;
735 parameter SUB = 6'b100010;
736 parameter SUBU = 6'b100011;
737 parameter AND = 6'b100100;
738 parameter OR = 6'b100101;
739 parameter XOR = 6'b100110;
740 parameter NOR = 6'b100111;
741 parameter SLT = 6'b101010;
742 parameter SLTU = 6'b101011;
743 parameter SLL = 6'b000000;
744 parameter SRL = 6'b000010;
745 parameter SRA = 6'b000011;
746 parameter SLLV = 6'b000100;
747 parameter SRLV = 6'b000110;
748 parameter SRAV = 6'b000111;
749 parameter JR = 6'b001000;
750 parameter LUI = 6'b001111;
751
752 wire signed [31:0] a_signed;
753 wire signed [31:0] b_signed;
754
755 reg [32:0] res;
756
757 assign a_signed = a;
758 assign b_signed = b;
759 assign r = res[31:0];
760
761 assign flag = (aluc == SLT || aluc == SLTU) ? ((aluc == SLT) ? (a_signed < b_signed) : (a < b)) : 1'bz;
762 assign zero = (res == 32'b0) ? 1'b1 : 1'b0;
763
764 always @ (a or b or aluc)
765 begin
766     case(aluc)
767         ADD: begin
768             res <= a_signed + b_signed;
769         end
770         ADDU: begin
771             res <= a + b;
772         end
773         SUB: begin

```

```

756         res <= a_signed - b_signed;
757     end
758     SUBU: begin
759         res <= a - b;
760     end
761     AND: begin
762         res <= a & b;
763     end
764     OR: begin
765         res <= a | b;
766     end
767     XOR: begin
768         res <= a ^ b;
769     end
770     NOR: begin
771         res <= ~(a | b);
772     end
773     SLT: begin
774         res <= a_signed < b_signed ? 1 : 0;
775     end
776     SLTU: begin
777         res <= a < b ? 1 : 0;
778     end
779     SLL: begin
780         res <= b << a;
781     end
782     SRL: begin
783         res <= b >> a;
784     end
785     SRA: begin
786         res <= b_signed >>> a_signed;
787     end
788     SLLV: begin
789         res <= b << a[4:0];
790     end
791     SRLV: begin
792         res <= b >> a[4:0];
793     end
794     SRAV: begin
795         res <= b_signed >>> a_signed[4:0];
796     end
797     LUI: begin
798         res <= {a[15:0], 16'h0000};
799     end
800     default:
801     begin
802         res <= 32'bz;
803     end
804 endcase
805 end
806 endmodule

```

Our expert-written area and power optimized solution: Shared adder for arithmetic, simplified flag logic and operand reuse, leading to 26% area reduction. Operand gating and early zeroing for large shifts to cut switching activity, for 4% power reduction.

```

795     1
796     2 wire sub_mode = (aluc==SUB) | (aluc==SUBU) | (aluc==SLT) | (aluc==SLTU);
797     3 wire [31:0] b_eff = sub_mode ? ~b : b;
798     4 wire cin = sub_mode;
799     5 wire [32:0] sum33 = {1'b0,a} + {1'b0,b_eff} + cin;
800     6
801     7 wire [31:0] add_res = sum33[31:0];
802     8 wire add_carry = sum33[32];
803     9
804     10 wire ovf = (a[31]^add_res[31]) & (b_eff[31]^add_res[31]);
805     11
806     12 wire signed_lt = add_res[31] ^ ovf;
807     13 wire uns_lt = ~add_carry;
808     14 wire [31:0] slt_res = {31'b0, signed_lt};
809     15 wire [31:0] sltu_res = {31'b0, uns_lt};
810     16
811     17 wire [31:0] and_res = a & b;
812     18 wire [31:0] or_res = a | b;
813     19 wire [31:0] xor_res = a ^ b;
814     20 wire [31:0] nor_res = ~(a | b);
815     21
816     22 wire [4:0] sa5 = a[4:0];
817     23 wire any_hi = |a[31:5];
818     24

```

```

810 25     wire [31:0] sll_full = any_hi ? 32'b0      : (b << sa5);
811 26     wire [31:0] srl_full = any_hi ? 32'b0      : (b >> sa5);
812 27     wire [31:0] sra_full = any_hi ? {32{b[31]}} : ($signed(b) >>> sa5);
813 28
814 29     wire [31:0] sllv_res = (b << a[4:0]);
815 30     wire [31:0] srlv_res = (b >> a[4:0]);
816 31     wire [31:0] srav_res = ($signed(b) >>> a[4:0]);
817 32
818 33     wire [31:0] lui_res = {a[15:0], 16'h0000};
819 34
820 35     reg [31:0] r_int;
821 36     always @* begin : result_mux
822 37         (* parallel_case, full_case *)
823 38         case (aluc)
824 39             ADD, ADDU: r_int = add_res;
825 40             SUB, SUBU: r_int = add_res;
826 41             AND:      r_int = and_res;
827 42             OR:       r_int = or_res;
828 43             XOR:      r_int = xor_res;
829 44             NOR:      r_int = nor_res;
830 45             SLT:      r_int = slt_res;
831 46             SLTU:     r_int = sltu_res;
832 47             SLL:      r_int = sll_full;
833 48             SRL:      r_int = srl_full;
834 49             SRA:      r_int = sra_full;
835 50             SLLV:     r_int = sllv_res;
836 51             SRLV:     r_int = srlv_res;
837 52             SRAV:     r_int = srav_res;
838 53             LUI:      r_int = lui_res;
839 54             JR:       r_int = 32'bz;
840 55             default:  r_int = 32'bz;
841 56         endcase
842 57     end
843 58
844 59     assign r      = r_int;
845 60     assign zero   = ~(r_int);
846 61
847 62     assign carry   = 1'bz;
848 63     assign overflow = 1'bz;
849 64     assign negative = 1'bz;
850 65     assign flag    = (aluc==SLT) ? signed_lt :
851 66                     (aluc==SLTU) ? uns_lt   :
852 67                     1'bz;
853 68 endmodule

```

Our expert-written delay optimized solution: Parallel datapaths with one-hot muxing for shallow critical path, leading to 26% delay reduction.

```

843 1     wire signed [31:0] a_signed = a;
844 2     wire signed [31:0] b_signed = b;
845 3
846 4     wire [31:0] add_u  = a + b;
847 5     wire [31:0] sub_u  = a - b;
848 6     wire signed [31:0] add_s = a_signed + b_signed;
849 7     wire signed [31:0] sub_s = a_signed - b_signed;
850 8
851 9     wire [31:0] and_res = a & b;
852 10    wire [31:0] or_res  = a | b;
853 11    wire [31:0] xor_res = a ^ b;
854 12    wire [31:0] nor_res = ~(a | b);
855 13
856 14    wire slt_res  = (a_signed < b_signed);
857 15    wire sltu_res = (a < b);
858 16
859 17    wire [4:0] shamt5 = a[4:0];
860 18
861 19
862 20    wire [31:0] sll_full  = (b << a);           // full 'a'
863 21    wire [31:0] srl_full  = (b >> a);           // full 'a'
864 22    wire [31:0] sra_full  = ($signed(b) >>> a_signed); // full signed 'a'
865 23    wire [31:0] sllv_res  = (b << shamt5);
866 24    wire [31:0] srlv_res  = (b >> shamt5);
867 25    wire [31:0] srav_res  = ($signed(b) >>> shamt5);
868 26
869 27    wire [31:0] lui_res = {a[15:0], 16'h0000};
870 28
871 29    wire sel_ADD  = (aluc==ADD);
872 30    wire sel_ADDU = (aluc==ADDU);
873 31    wire sel_SUB  = (aluc==SUB);
874 32    wire sel_SUBU = (aluc==SUBU);

```



```

864 33     wire sel_AND  = (aluc==AND);
865 34     wire sel_OR   = (aluc==OR);
866 35     wire sel_XOR   = (aluc==XOR);
867 36     wire sel_NOR   = (aluc==NOR);
868 37     wire sel_SLT   = (aluc==SLT);
869 38     wire sel_SLTU  = (aluc==SLTU);
870 39     wire sel_SLL   = (aluc==SLL);
871 40     wire sel_SRL   = (aluc==SRL);
872 41     wire sel_SRA   = (aluc==SRA);
873 42     wire sel_SLLV  = (aluc==SLLV);
874 43     wire sel_SRLV  = (aluc==SRLV);
875 44     wire sel_SRAV  = (aluc==SRAV);
876 45     wire sel_LUI   = (aluc==LUI);
877 46     wire sel_JR    = (aluc==JR);
878 47
879 48     wire any_sel = sel_ADD|sel_ADDU|sel_SUB|sel_SUBU|sel_AND|sel_OR|sel_XOR|sel_NOR|
880 49                  sel_SLT|sel_SLTU|sel_SLL|sel_SRL|sel_SRA|sel_SLLV|sel_SRLV|sel_SRAV|sel_LUI;
881 50
882 51     wire [31:0] r_known =
883 52         (sel_ADD ? add_s : 32'b0) |
884 53         (sel_ADDU ? add_u : 32'b0) |
885 54         (sel_SUB ? sub_s : 32'b0) |
886 55         (sel_SUBU ? sub_u : 32'b0) |
887 56         (sel_AND ? and_res : 32'b0) |
888 57         (sel_OR ? or_res : 32'b0) |
889 58         (sel_XOR ? xor_res : 32'b0) |
890 59         (sel_NOR ? nor_res : 32'b0) |
891 60         (sel_SLT ? {31'b0, slt_res} : 32'b0) |
892 61         (sel_SLTU ? {31'b0, sltu_res} : 32'b0) |
893 62         (sel_SLL ? sll_full : 32'b0) |
894 63         (sel_SRL ? srl_full : 32'b0) |
895 64         (sel_SRA ? sra_full : 32'b0) |
896 65         (sel_SLLV ? sllv_res : 32'b0) |
897 66         (sel_SRLV ? srlv_res : 32'b0) |
898 67         (sel_SRAV ? srav_res : 32'b0) |
899 68         (sel_LUI ? lui_res : 32'b0);
900 69
901 70     assign r = (any_sel && !sel_JR) ? r_known : 32'bz;
902 71
903 72     assign zero = (r == 32'b0) ? 1'b1 : 1'b0;
904 73
905 74     assign flag = (sel_SLT) ? slt_res :
906 75                  (sel_SLTU) ? sltu_res :
907 76                  1'bz;
908 77
909 78     assign carry = 1'bz;
910 79     assign negative = 1'bz;
911 80     assign overflow = 1'bz;
912 81 endmodule

```

A.3.2 PROBLEM #68: RTLLM FSM

Problem description: Implement a finit state machine (FSM) that detects the input sequence 10011 on a single-bit input stream. The module has three inputs: the serial input bit (IN), the clock (CLK), and a synchronous reset (RST). It produces one output, MATCH, which is asserted high when the specified sequence is recognized. The FSM supports continuous input and loop detection. When reset is active, the FSM initializes and MATCH is cleared to 0. The output MATCH is asserted during the cycle when the last 1 of the target sequence is received, and the design ensures that repeated or overlapping patterns (e.g., 100110011) correctly generate multiple match pulses.

Benchmark solution: States are binary-encoded, with sequential next-state logic in a Mealy FSM while output occupies a register.

```

909 1 module unopt_model (
910 2     input wire IN,
911 3     input wire CLK,
912 4     input wire RST,
913 5     output wire MATCH
914 6 );
915 7
916 8 reg [2:0] ST_cr, ST_nt;
917 9
918 10 parameter s0 = 3'b000;
919 11 parameter s1 = 3'b001;
920 12 parameter s2 = 3'b010;
921 13 parameter s3 = 3'b011;
922 14 parameter s4 = 3'b100;

```

```

918 15 parameter s5 = 3'b101;
919 16
920 17 always@(posedge CLK or posedge RST) begin
921 18     if(RST)
922 19         ST_cr <= s0;
923 20     else
924 21         ST_cr <= ST_nt;
925 22 end
926 23
927 24 always@(*) begin
928 25     case(ST_cr)
929 26         s0:begin
930 27             if (IN==0)
931 28                 ST_nt = s0;
932 29             else
933 30                 ST_nt = s1;
934 31         end
935 32         s1:begin
936 33             if (IN==0)
937 34                 ST_nt = s2;
938 35             else
939 36                 ST_nt = s1;
940 37         end
941 38         s2:begin
942 39             if (IN==0)
943 40                 ST_nt = s3;
944 41             else
945 42                 ST_nt = s1;
946 43         end
947 44         s3:begin
948 45             if (IN==0)
949 46                 ST_nt = s0;
950 47             else
951 48                 ST_nt = s4;
952 49         end
953 50         s4:begin
954 51             if (IN==0)
955 52                 ST_nt = s2;
956 53             else
957 54                 ST_nt = s5;
958 55         end
959 56         s5:begin
960 57             if (IN==0)
961 58                 ST_nt = s2;
962 59             else
963 60                 ST_nt = s1;
964 61         end
965 62     endcase
966 63 end
967 64
968 65 always@(*) begin
969 66     if(RST)
970 67         MATCH <= 0;
971 68     else if (ST_cr == s4 && IN == 1)
972 69         MATCH <= 1;
973 70     else
974 71         MATCH <= 0;
975 72 end
976 73
977 74 endmodule
978 75
979 76
980 77
981 78
982 79
983 80

```

Our expert-written area and power optimized solution: Casez-based transitions and direct Mealy output computation, removing extra register, leading to 32% area reduction. Compact binary encoding and casez-based transitions to cut toggling for 23% power reduction.

```

966 1 localparam [2:0] s0=3'b000, s1=3'b001, s2=3'b010,
967 2         s3=3'b011, s4=3'b100, s5=3'b101;
968 3
969 4 reg [2:0] ST_cr, ST_nt;
970 5
971 6 always @(posedge CLK or posedge RST) begin
972 7     if (RST)
973 8         ST_cr <= s0;
974 9     else

```

```

972     ST_cr <= ST_nt;
973 end
974 always @* begin
975     ST_nt = s0;
976     casez ({ST_cr, IN})
977         // s0: 0->s0, 1->s1
978         {s0,1'b0}: ST_nt = s0;
979         {s0,1'b1}: ST_nt = s1;
980
981         // s1: 0->s2, 1->s1
982         {s1,1'b0}: ST_nt = s2;
983         {s1,1'b1}: ST_nt = s1;
984
985         // s2: 0->s3, 1->s1
986         {s2,1'b0}: ST_nt = s3;
987         {s2,1'b1}: ST_nt = s1;
988
989         // s3: 0->s0, 1->s4
990         {s3,1'b0}: ST_nt = s0;
991         {s3,1'b1}: ST_nt = s4;
992
993         // s4: 0->s2, 1->s5
994         {s4,1'b0}: ST_nt = s2;
995         {s4,1'b1}: ST_nt = s5;
996
997         // s5: 0->s2, 1->s1
998         {s5,1'b0}: ST_nt = s2;
999         {s5,1'b1}: ST_nt = s1;
1000     default: ST_nt = s0;
1001 endcase
1002 end
1003 assign MATCH = (ST_cr == s4) & IN;
1004 endmodule

```

Our expert-written delay optimized solution: One-hot state encoding with pre-decoded inputs and Moore-style output, leading to 23% delay reduction.

```

1000 1
1001 2 reg [5:0] S, S_next;
1002 3
1003 4 reg [5:0] S, S_next;
1004 5
1005 6 always @(posedge CLK or posedge RST) begin
1006 7     if (RST)
1007 8         S <= 6'b000001; // s0
1008 9     else
1009 10         S <= S_next;
1010 11 end
1011 12
1012 13 wire in1 = IN;
1013 14 wire in0 = ~IN;
1014 15
1015 16 always @* begin
1016 17     S_next[0] = (S[0] & in0) | (S[3] & in0); // -> s0
1017 18     S_next[1] = (S[0] & in1) | (S[1] & in1) | (S[2] & in1) // -> s1
1018 19     | (S[5] & in1);
1019 20     S_next[2] = (S[1] & in0) | (S[4] & in0) | (S[5] & in0); // -> s2
1020 21     S_next[3] = (S[2] & in0); // -> s3
1021 22     S_next[4] = (S[3] & in1); // -> s4
1022 23     S_next[5] = (S[4] & in1); // -> s5
1023 24 end
1024 25
1025 26 always @(posedge CLK or posedge RST) begin
1026 27     if (RST)
1027 28         MATCH <= 1'b0;
1028 29     else
1029 30         MATCH <= S[5];
1030 31 end
1031 32 endmodule

```

A.3.3 PROBLEM #87: VERILOGEVAL PROB095

Problem description: Implement a module that generates a control signal (shift_ena) for a shift register. The module has a clock input (clk), a synchronous active-high reset (reset), and a single

output (shift_ena). The functionality requires that when the FSM is reset, the shift_ena signal is asserted high for exactly four consecutive clock cycles before being deasserted permanently. After this sequence, shift_ena remains low indefinitely until another reset occurs, at which point the behavior repeats. All sequential operations are triggered on the positive edge of the clock.

Benchmark solution: Uses explicit states with next-state logic to drive shift_ena.

```

1031 1 module unopt_model (
1032 2     input  clk,
1033 3     input  reset,
1034 4     output reg shift_ena
1035 5 );
1036 6     parameter B0=0, B1=1, B2=2, B3=3, Done=4;
1037 7
1038 8     reg [2:0] state, next;
1039 9
1040 10    always @* begin
1041 11        case (state)
1042 12            B0: next = B1;
1043 13            B1: next = B2;
1044 14            B2: next = B3;
1045 15            B3: next = Done;
1046 16            Done: next = Done;
1047 17            default: next = B0;
1048 18        endcase
1049 19    end
1050 20
1051 21    always @(posedge clk) begin
1052 22        if (reset) begin
1053 23            state <= B0;
1054 24            shift_ena <= 1'b1;
1055 25        end else begin
1056 26            state <= next;
1057 27            shift_ena <= (next != Done);
1058 28        end
1059 29    end
1060 30 endmodule

```

Our expert-written area and power optimized solution: Minimized register width, using a 2-bit counter, and compact comparator logic, leading to 47% area reduction. Small counter reused which reduced toggling activity to minimize dynamic power, for 46% power reduction.

```

1061 1     reg [1:0] counter; // 2 bits are enough to count up to 4
1062 2
1063 3     always @(posedge clk) begin
1064 4         if (reset) begin
1065 5             counter <= 2'b00; // Reset counter
1066 6             shift_ena <= 1'b1; // Enable on reset
1067 7         end else if (counter < 2'b11) begin
1068 8             counter <= counter + 1; // Increment counter
1069 9             shift_ena <= 1'b1; // Keep shift_ena high while counting
1070 10        end else begin
1071 11            shift_ena <= 1'b0; // Disable after 4 cycles
1072 12        end
1073 13    end
1074 14 endmodule

```

Our expert-written delay optimized solution: Wider counter, using 3 bits, to simplify comparison and reduce logic depth on the critical path, leading to 37% delay reduction.

```

1075 1     reg [2:0] count; // 3-bit counter to count 4 cycles
1076 2
1077 3     always @(posedge clk) begin
1078 4         if (reset) begin
1079 5             count <= 3'b000; // Reset count
1080 6             shift_ena <= 1'b1; // Enable shift initially
1081 7         end else if (count < 3'b011) begin
1082 8             count <= count + 1; // Increment count
1083 9             shift_ena <= 1'b1; // Keep shift enabled
1084 10        end else begin
1085 11            shift_ena <= 1'b0; // Disable shift after 4 cycles
1086 12        end
1087 13    end
1088 14 endmodule

```

A.3.4 PROBLEM #104: VERILOG EVAL PROB136

Problem description: Implement a cellular automaton game, similar to Conway’s Game of Life, on a 16x16 grid. The grid is represented as a 256-bit vector (q), where each row of 16 cells maps to a sub-vector, and each cell can be alive (1) or dead (0). The module has a clock input (clk), a load signal (load) for synchronously loading an initial 256-bit state (data) into q, and produces the updated 256-bit grid state as output. At every positive clock edge, the grid advances by one timestep, with each cell’s next state determined by its number of neighbors: cells die with fewer than 2 or more than 3 neighbors, remain unchanged with exactly 2 neighbors, and become alive with exactly 3 neighbors. The grid is modeled as a toroid, meaning edges wrap around so that cells on the boundaries consider neighbors from the opposite side.

Benchmark solution: Straightforward RTL with per-cell neighbor recomputation and sequential summing.

```

1092 module unopt_model (
1093     input clk,
1094     input load,
1095     input [255:0] data,
1096     output reg [255:0] q
1097 );
1098
1099     logic [323:0] q_pad;
1100     always@(*) begin
1101         for (int i=0; i<16; i++)
1102             q_pad[18*(i+1)+1 +: 16] = q[16*i +: 16];
1103         q_pad[1 +: 16] = q[16*15 +: 16];
1104         q_pad[18*17+1 +: 16] = q[0 +: 16];
1105
1106         for (int i=0; i<18; i++) begin
1107             q_pad[i*18] = q_pad[i*18+16];
1108             q_pad[i*18+17] = q_pad[i*18+1];
1109         end
1110     end
1111
1112     always @(posedge clk) begin
1113         for (int i=0; i<16; i++)
1114             for (int j=0; j<16; j++) begin
1115                 q[i*16+j] <=
1116                     ((q_pad[(i+1)*18+j+1 -1+18] + q_pad[(i+1)*18+j+1 +18] + q_pad[(i+1)*18+j+1 +1+18] +
1117                      q_pad[(i+1)*18+j+1 -1] + q_pad[(i+1)*18+j+1+1] +
1118                      q_pad[(i+1)*18+j+1 -1-18] + q_pad[(i+1)*18+j+1 -18] + q_pad[(i+1)*18+j+1 +1-18]) & 3'h7 | q[i*16+j]) == 3'h3;
1119             end
1120         if (load)
1121             q <= data;
1122     end
1123 endmodule

```

Our expert-written area and power optimized solution: Sharing per-row horizontal sums and bitwise rotations for toroidal wrap, minimizing summations for 19% area reduction. Computation reuse decreasing toggling, leading to 36% power reduction.

```

1120 // --- Helpers -----
1121 function automatic [7:0] idx(input [3:0] r, input [3:0] c);
1122     idx = {r, c}; // r*16 + c
1123 endfunction
1124
1125 // Bit-rotate wires (toroidal wrap) - wiring only (no logic area)
1126 function automatic [15:0] roll(input [15:0] x); roll = {x[14:0], x[15]}; endfunction
1127 function automatic [15:0] rorl(input [15:0] x); rorl = {x[0], x[15:1]}; endfunction
1128
1129 // Add-three 1-bit vectors in parallel: returns {carry,sum}
1130 // a+b+c = sum ^ (2*carry) per bit
1131 function automatic [31:0] add3_vec(input [15:0] a, input [15:0] b, input [15:0] c);
1132     add3_vec[15:0] = a ^ b ^ c; // sum (LSB)
1133     add3_vec[31:16] = (a & b) | (a & c) | (b & c); // carry (means +2)
1134 endfunction
1135
1136 // --- Unpack rows (wires) -----
1137 wire [15:0] row [15:0];
1138 genvar ur;
1139 generate
1140     for (ur = 0; ur < 16; ur = ur + 1) begin : UNPACK
1141         assign row[ur] = q[{ur[3:0], 4'b0000} +: 16];
1142     end
1143 endgenerate

```

```

1134   end
1135   endgenerate
1136
1137   // --- Precompute per-row horizontal neighbors (shared) -----
1138   wire [15:0] rol [15:0], ror [15:0];
1139   wire [15:0] sTrip [15:0], cTrip [15:0]; // for (L,C,R) of each row (0..3)
1140   wire [15:0] sPair [15:0], cPair [15:0]; // for (L,R) of each row (0..2)
1141
1142   genvar hr;
1143   generate
1144   for (hr = 0; hr < 16; hr = hr + 1) begin : HROW
1145     assign rol[hr] = roll(row[hr]);
1146     assign ror[hr] = rorl(row[hr]);
1147
1148     // Triplet = left + center + right (encoded as s + 2*c)
1149     wire [31:0] trip_pack = add3_vec(rol[hr], row[hr], ror[hr]);
1150     assign sTrip[hr] = trip_pack[15:0];
1151     assign cTrip[hr] = trip_pack[31:16];
1152
1153     // Pair = left + right
1154     assign sPair[hr] = rol[hr] ^ ror[hr];
1155     assign cPair[hr] = rol[hr] & ror[hr];
1156   end
1157   endgenerate
1158
1159   integer r;
1160   reg [255:0] nxt;
1161
1162   reg [3:0] rn, rp;
1163   reg [15:0] sT, cT, sM, cM, sB, cB;
1164   reg [15:0] sS, cS; // sum of (sT + sM + sB) as s + 2*c
1165   reg [15:0] U_is0, U_ge2, U_is1; // onehot(U) for U = cS + cT + cM + cB
1166   reg [15:0] is3, is2;
1167
1168   always @* begin
1169     nxt = '0;
1170
1171     for (r = 0; r < 16; r = r + 1) begin
1172       rn = (r == 0) ? 4'd15 : r - 1;
1173       rp = (r == 15) ? 4'd0 : r + 1;
1174
1175       sT = sTrip[rn]; cT = cTrip[rn]; // top triplet from row r-1
1176       sM = sPair[r ]; cM = cPair[r ]; // middle pair (no center) from row r
1177       sB = sTrip[rp]; cB = cTrip[rp]; // bottom triplet from row r+1
1178
1179       {cS, sS} = add3_vec(sT, sM, sB);
1180
1181       U_is0 = ~(cS | cT | cM | cB);
1182       U_ge2 = ( (cS & cT) | (cS & cM) | (cS & cB)
1183                | (cT & cM) | (cT & cB) | (cM & cB) );
1184       U_is1 = ~(U_is0 | U_ge2);
1185
1186       is3 = sS & U_is1;
1187       is2 = ~sS & U_is1;
1188
1189       nxt[{r[3:0], 4'b0000} +: 16] = is3 | (row[r] & is2);
1190     end
1191   end
1192
1193   always @(posedge clk) begin
1194     if (load)
1195       q <= data;
1196     else
1197       q <= nxt;
1198   end
1199 endmodule

```

Our expert-written delay optimized solution: Parallel neighbor summation with carry-save adder tree and direct decode for 2 and 3, for shallow critical path, leading to 37% delay reduction.

```

1200 function automatic [7:0] idx(input [3:0] r, input [3:0] c);
1201   idx = {r, c};
1202 endfunction
1203
1204 function automatic [15:0] roll(input [15:0] x); roll = {x[14:0], x[15]}; endfunction
1205 function automatic [15:0] rorl(input [15:0] x); rorl = {x[0], x[15:1]}; endfunction
1206
1207 function automatic [31:0] add3_vec(input [15:0] a, input [15:0] b, input [15:0] c);
1208   add3_vec[15:0] = a ^ b ^ c; // s
1209   add3_vec[31:16] = (a & b) | (a & c) | (b & c); // c (>=2)

```

```

1188 11 endfunction
1189 12
1190 13 integer r;
1191 14 reg [255:0] nxt;
1192 15
1193 16 reg [3:0] rn, rp;
1194 17 reg [15:0] ru, r0, rd;
1195 18 reg [15:0] ru_l, ru_c, ru_r;
1196 19 reg [15:0] r0_l, r0_r;
1197 20 reg [15:0] rd_l, rd_c, rd_r;
1198 21
1199 22 reg [15:0] sT, cT, sM, cM, sB, cB, sS, cS;
1200 23 reg [15:0] U_is0, U_ge2, U_is1; // onehot decode for U = cS+cT+cM+cB
1201 24 reg [15:0] is3, is2; // neighbor count ==3 / ==2
1202 25
1203 26 always @* begin
1204 27     nxt = '0;
1205 28
1206 29     for (r = 0; r < 16; r = r + 1) begin
1207 30         rn = (r == 0) ? 4'd15 : r - 1;
1208 31         rp = (r == 15) ? 4'd0 : r + 1;
1209 32
1210 33         ru = q[{rn, 4'b0000} +: 16];
1211 34         r0 = q[{r, 4'b0000} +: 16];
1212 35         rd = q[{rp, 4'b0000} +: 16];
1213 36
1214 37         ru_l = roll(ru); ru_c = ru; ru_r = rorl(ru);
1215 38         r0_l = roll(r0); r0_r = rorl(r0);
1216 39         rd_l = roll(rd); rd_c = rd; rd_r = rorl(rd);
1217 40
1218 41         {cT, sT} = add3_vec(ru_l, ru_c, ru_r); // counts 0..3
1219 42         sM = (r0_l ^ r0_r); // pair: 0..2
1220 43         cM = (r0_l & r0_r);
1221 44         {cB, sB} = add3_vec(rd_l, rd_c, rd_r);
1222 45
1223 46         {cS, sS} = add3_vec(sT, sM, sB);
1224 47
1225 48         U_is0 = ~(cS | cT | cM | cB);
1226 49         U_ge2 = ( (cS & cT) | (cS & cM) | (cS & cB)
1227 50                 | (cT & cM) | (cT & cB) | (cM & cB) );
1228 51         U_is1 = ~(U_is0 | U_ge2);
1229 52
1230 53         is3 = sS & U_is1;
1231 54         is2 = ~sS & U_is1;
1232 55
1233 56         nxt[{r[3:0], 4'b0000} +: 16] = is3 | (r0 & is2);
1234 57     end
1235 58 end
1236 59
1237 60 always @(posedge clk) begin
1238 61     if (load)
1239 62         q <= data;
1240 63     else
1241 64         q <= nxt;
1242 65     end
1243 66 endmodule

```

A.4 PASS@K DEFINITION

The pass@k metric, defined in equation 4, is used for measuring the functional correctness of LLM-generated code. Here, N is the total number of problems in the evaluation set. For a given problem i , we generate n_i samples and evaluate their functional correctness, obtaining c_i correct samples. The metric estimates the probability that at least one of the correct solution appears when drawing k samples.

$$\text{pass}@k = \mathbb{E}_{i=1}^N \left[1 - \frac{C(n_i - c_i, k)}{C(n_i, k)} \right] \quad (4)$$

where $C(n, k)$ denotes the binomial coefficient 'n choose k' and the term $1 - \frac{C(n_i - c_i, k)}{C(n_i, k)}$ represents the probability that at least one of the k samples is correct from the n_i total samples.

The metric averages this probability across all N problems to provide an unbiased estimate of the model's functional correctness when generating k samples per problem.

Implementation Note: Following standard practice, we generate $n \geq k$ samples per problem (typically $n = 10$) and compute pass@k for various values of k (e.g., $k \in 1, 5, 10$) to obtain unbiased estimates without requiring k separate evaluation runs.

A.5 SAMPLING DIVERSITY

In Figure 7, we show three LLM-generated Verilog modules of a parallel-in-serial-out shift register, targeting area optimization, with normalized area efficiency values of 0.924, 0.963 and 1.0, respectively. This diversity demonstrates why eff@k with $k > 1$ is valuable: generating multiple samples increases the likelihood of obtaining high-quality implementations. The optimal sample (rightmost) avoids unnecessary counters and state tracking present in the other solutions, reducing both area and complexity, a correct solution that might not appear in a single generation attempt but becomes accessible when drawing multiple samples.

```

module opt_model #(parameter DATA_WIDTH = 32) (
    input clk,
    input resetn,
    input [DATA_WIDTH-1:0] din,
    input din_en,
    output logic dout
);
    reg [DATA_WIDTH-1:0] shift_reg;
    reg [$clog2(DATA_WIDTH)-1:0] bit_index;
    always @(posedge clk or negedge resetn) begin
        if (!resetn) begin
            shift_reg <= {DATA_WIDTH{1'b0}};
            bit_index <= 0;
        end else if (din_en) begin
            shift_reg <= din;
            bit_index <= 0;
            dout <= din[0];
        end else if (bit_index < DATA_WIDTH-1) begin
            dout <= shift_reg[bit_index + 1];
            bit_index <= bit_index + 1;
        end else begin
            dout <= 1'b0;
        end
    end
endmodule

```

```

module opt_model #(parameter DATA_WIDTH = 32) (
    input clk,
    input resetn,
    input [DATA_WIDTH-1:0] din,
    input din_en,
    output logic dout
);
    reg [DATA_WIDTH-1:0] shift_reg;
    reg [$clog2(DATA_WIDTH)-1:0] bit_index;
    always @(posedge clk) begin
        if (din_en) begin
            shift_reg <= din;
            bit_index <= 0;
            dout <= din[0];
        end else if (bit_index < DATA_WIDTH-1) begin
            dout <= shift_reg[bit_index + 1];
            bit_index <= bit_index + 1;
        end else begin
            dout <= 1'b0;
        end
    end
endmodule

```

```

module opt_model #(parameter DATA_WIDTH = 32) (
    input clk,
    input resetn,
    input [DATA_WIDTH-1:0] din,
    input din_en,
    output logic dout
);
    logic [DATA_WIDTH-1:0] temp;
    always @(posedge clk or negedge resetn) begin
        if (!resetn) begin
            temp <= 0;
        end else if (din_en) begin
            temp <= din;
        end else begin
            temp <= temp >> 1;
        end
    end
    assign dout = temp[0];
endmodule

```

(a) Top generation at $k = 1$

(b) Top generation at $k = 5$

(c) Top generation at $k = 10$

Figure 7: Three area-optimized implementations of the `piso_shift_register` module (Problem #16) generated at $k \in \{1, 5, 10\}$. The circuit shifts the least significant bit of a multi-bit input `din` to the single-bit output `dout` sequentially, starting when `din_en` goes high. All designs are functionally correct but structurally diverse.

A.6 FAILURE ANALYSIS

In Figure. 8, we visualize the set of problems that have high pass@k and low eff@k to get better insight on the set of problems that are hard to optimize.

A.7 PROMPTING WITH OPTIMIZATION STRATEGIES

To assess whether the relatively lower eff@k scores stem from missing optimization guidance, we repeated the experiment using the first prompting strategy, where the LLM rewrites an unoptimized Verilog module into a more efficient implementation, but this time we explicitly provided the optimization goals and permissible strategies for each metric. The results, summarized in Table 5, show that explicit guidance yields only marginal and model-dependent improvements. A few models (e.g., Mixtral, StarCoder2-15B, CodeLlama-70B, and CodeV-QC-7B) show moderate gains across several efficiency metrics, while others improve only in isolated cases (e.g., GPT-4o-mini improves mainly in power, Qwen2.5-Coder-7B only in delay). However, for several strong models, including DeepSeek-Chat and Llama-3.3-70B, efficiency scores worsen despite being given optimization strategies. Overall, the gap between pass@k and eff@k remains substantial, indicating that the limited efficiency performance is not primarily caused by insufficient prompt specification, but rather reflects the underlying difficulty of generating functionally correct and resource-efficient RTL simultaneously.

Table 5: Efficiency results when using the first prompting strategy with explicit optimization goals and metric-specific optimization strategies. While a few models exhibit modest or selective improvements, most models show limited or no gains, and several regress. Overall, providing optimization guidance does not substantially narrow the gap between pass@k and eff@k.

Model	pass@1	pass@5	pass@10	eff@1			eff@5			eff@10		
				Area	Delay	Power	Area	Delay	Power	Area	Delay	Power
GPT-4o-mini	0.471	0.667	0.725	0.435	0.440	0.423	0.622	0.637	0.642	0.696	0.693	0.716
DeepSeek-Chat	0.502	0.728	0.795	0.487	0.496	0.485	0.717	0.719	0.709	0.795	0.789	0.774
Llama-3.3-70B-Instruct	0.429	0.666	0.752	0.400	0.412	0.404	0.626	0.666	0.627	0.703	0.753	0.714
Mixtral-8x7B-v0.1	0.271	0.562	0.696	0.232	0.255	0.244	0.491	0.528	0.508	0.607	0.648	0.632
starcode2-15b-instruct-v0.1	0.686	0.969	0.985	0.632	0.653	0.615	0.893	0.917	0.888	0.904	0.935	0.904
CodeLlama-70b-Instruct-hf	0.599	0.904	0.944	0.538	0.559	0.547	0.834	0.857	0.846	0.877	0.905	0.895
DeepSeek-Coder-33B	0.6088	0.9358	0.9766	0.5539	0.5706	0.5648	0.8745	0.8948	0.887	0.9270	0.9514	0.9191
Qwen2.5-Coder-7B-Inst	0.443	0.756	0.839	0.399	0.421	0.389	0.675	0.735	0.692	0.745	0.834	0.785
yang-z/CodeV-QC-7B	0.396	0.662	0.725	0.360	0.359	0.347	0.593	0.638	0.601	0.644	0.718	0.667
RTLCoder-DeepSeek-V1	0.470	0.816	0.895	0.418	0.428	0.413	0.738	0.780	0.726	0.828	0.860	0.803
VeriThoughts-Inst.-7B	0.545	0.770	0.836	0.475	0.508	0.492	0.673	0.715	0.699	0.748	0.774	0.758

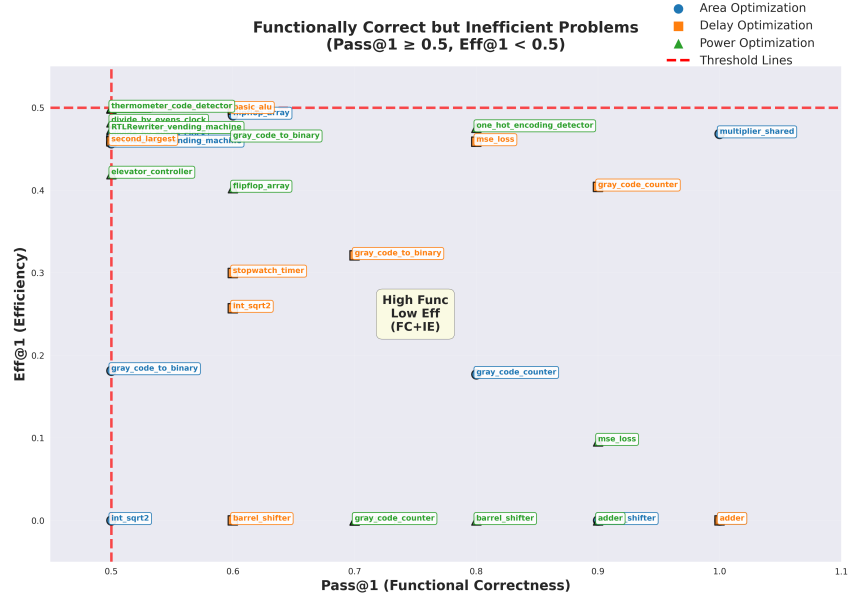


Figure 8: Quadrant plot showing the correlation between functional correctness ($Pass@1$) and synthesis efficiency ($Eff@1$) across area, delay, and power objectives.