

AI for VLSI

A Practical Guide to AI-Driven RTL Design and Verification



P.L. Ho

First Edition · 2026

Table of Contents

Introduction

AI/ML Primer for VLSI Engineers

LLMs for RTL Generation

AI-Assisted RTL Code Review and Bug Detection

AI for Functional Verification

AI for Formal Verification

Tools, Frameworks, and Workflows

Getting Started and Future Outlook

AI Development Environment for VLSI Engineers

Appendix A: Prompt Library

Appendix B: Review Checklists

Appendix C: Metrics Dashboard Template

Appendix D: Data Governance and IP Policy

Appendix E: AI Rules Reference for VLSI Projects

References and Citation Notes

HOW TO USE THIS BOOK

This guide is designed for practicing VLSI engineers and managers. Focus first on Chapters 1, 3, 5, and 8 if your goal is immediate adoption. Start with Chapter 9 and Appendix E if you want to set up your AI development environment right away. Use Chapters 2, 4, 6, and 7 as deeper technical and tooling references.

As-of note: Tool capabilities and benchmark numbers in this book reflect the landscape as of 2026 and should be validated against current releases before procurement or sign-off decisions.

Evidence labels: **Production Practice** = used in shipping flows; **Research Result** = published but may require internal validation; **Speculative Direction** = plausible roadmap, not production-ready.

CHAPTER 1

Introduction

The Complexity Wall

Modern systems-on-chip contain billions of transistors. The Apple M3 Ultra integrates over 184 billion; NVIDIA's Blackwell B200 exceeds 208 billion [1][2]. At the 3 nm and 2 nm nodes now entering production, a single advanced SoC may contain hundreds of IP blocks, dozens of clock domains, and millions of lines of RTL. Design teams routinely manage codebases of 5–10 million lines of SystemVerilog, VHDL, and C/C++ testbench code [3].

This complexity has grown exponentially for decades, roughly following Moore's Law. But engineering team sizes and EDA tool productivity have not kept pace. The result is what the semiconductor industry calls the **design productivity gap**—the widening chasm between what silicon technology can theoretically integrate and what human engineers can practically design and verify in a competitive product cycle.

KEY CONCEPT: THE DESIGN PRODUCTIVITY GAP

Transistor counts double approximately every two years, but engineering productivity (measured in transistors designed per engineer per day) improves at only 20–30% per year. This exponential divergence means that without new tools and methodologies, design teams would need to double in size every few years—an economic impossibility for all but the largest companies.

Verification is the dominant bottleneck. Industry surveys consistently report that 60–70% of total project effort goes to functional verification [4]. A complex SoC may require hundreds of millions of simulation cycles, thousands of formal properties, and months of emulation time before tapeout. Despite this effort, first-silicon bugs remain common—Synopsys reports that over 50% of designs require at least one respin, at a cost of \$5–50 million per iteration at advanced nodes [5].

How AI/ML Is Changing EDA

Artificial intelligence and machine learning are now making measurable inroads into nearly every stage of the chip design flow. This is not speculative futurism—production tools from all three major EDA vendors (Synopsys, Cadence, and Siemens EDA) already incorporate ML-driven optimization in synthesis, place-and-route, timing closure, and verification.

Production Practice: Many of the capabilities below are deployed today in commercial tools, though ROI depends on design class and integration maturity [6][7][8].

The applications span a wide spectrum:

- **RTL generation:** Large language models (LLMs) trained on hardware description languages can produce Verilog and SystemVerilog modules from natural-language specifications or partial code.
- **Code review and bug detection:** ML models trained on historical bug databases can flag likely defects—clock domain crossings, unintended latches, FSM deadlocks—before simulation.
- **Verification:** Reinforcement learning agents can steer constrained-random stimulus toward uncovered states, reducing time-to-coverage by 2–5×.
- **Physical design:** Google's AlphaChip (originally published as the Nature paper on chip floorplanning with RL) demonstrated superhuman macro placement quality. Commercial tools like Synopsys DSO.ai and Cadence Cerebrus now use RL and Bayesian optimization to explore synthesis and P&R parameter spaces autonomously.
- **Analog design:** ML-based sizing and layout optimization are reducing analog design cycles from weeks to hours.

From Rule-Based to Data-Driven

Traditional EDA tools operate on hand-crafted heuristics. A synthesis tool's optimization passes, a router's layer-assignment strategy, and a linter's rule set are all products of decades of expert engineering. These heuristics work well on average but struggle with the exponentially large search spaces of modern designs.

The shift to data-driven methods mirrors a transition that has already transformed software engineering, drug discovery, and materials science. In EDA, this transition has been gradual: early academic work in the 2010s applied random forests and

support vector machines to timing prediction and lithography hotspot detection. The transformer revolution of 2017–2023 opened the door to generative models for code and natural-language-driven design. By 2024–2025, fine-tuned LLMs for Verilog, graph neural networks for netlists, and RL agents for physical design had all moved from research prototypes to commercial or near-commercial deployment.

What This Guide Covers

This book is a practical reference for working VLSI engineers—RTL designers, verification engineers, physical design engineers, and engineering managers—who want to understand and adopt AI/ML in their workflows. It is not a machine learning textbook. We cover exactly as much ML theory as you need to make informed decisions about tools, evaluate vendor claims, and build internal capabilities.

Chapter 2 covers the essential ML concepts. Chapters 3–4 focus on the front-end design flow: RTL generation with LLMs and AI-assisted code review. Chapters 5–6 address functional and formal verification. Chapter 7 surveys commercial and open-source tools. Chapter 8 provides a practical roadmap for getting started. Chapter 9 covers the hands-on setup of AI development environments—project rules, the Model Context Protocol (MCP) for EDA tool integration, reusable AI skills, and advanced prompt engineering techniques.

Why Now?

Three trends have converged to make AI-driven chip design practical today:

1. **Compute:** GPU clusters and cloud infrastructure make it feasible to train and serve large models. A single NVIDIA H100 delivers more than 30× the training throughput of the V100s available just five years ago.
2. **Data:** Open-source HDL repositories (GitHub, OpenCores), academic benchmarks (IWLS, EPFL), and internal design databases provide the training corpora that ML models require. Companies like NVIDIA have demonstrated that fine-tuning on proprietary RTL corpora yields significantly better results than general-purpose models.
3. **Algorithms:** Transformer architectures, reinforcement learning from human feedback (RLHF), and retrieval-augmented generation (RAG) have dramatically improved the quality and controllability of generative AI systems.

The semiconductor industry is at an inflection point. AI will not replace chip designers—but chip designers who use AI will increasingly outperform those who do not.

CHAPTER 2

AI/ML Primer for VLSI Engineers

This chapter introduces the machine learning concepts you will encounter throughout the rest of the book. We deliberately frame everything in terms that electrical engineers already understand: transfer functions, optimization, signal processing, and state machines. If you have a solid grasp of linear algebra and probability from your undergraduate EE curriculum, you have all the mathematical background you need.

Neural Networks as Parameterized Transfer Functions

A neural network is, at its core, a parameterized nonlinear function $f(x; \theta)$ that maps an input vector x to an output vector y . Think of it the same way you think of a complex analog circuit: there is a set of parameters (weights and biases, analogous to resistor values and transistor sizes) that determine the input-output behavior. Training the network means finding the parameter values that minimize a cost function—exactly like optimizing a circuit's performance metrics.

A single neuron computes $y = \sigma(w \cdot x + b)$, where w is a weight vector, b is a bias, and σ is a nonlinear activation function (ReLU, sigmoid, or tanh). This is directly analogous to a weighted summer followed by a nonlinear transfer characteristic. Stack thousands of these neurons into layers, and you get a network capable of approximating arbitrarily complex functions—the Universal Approximation Theorem guarantees this, though it says nothing about how easy it is to find the right parameters.

KEY CONCEPT: TRAINING = OPTIMIZATION

Training a neural network is gradient-descent optimization of a loss function over a parameter space. If you have ever swept transistor widths to minimize power while meeting timing, you already understand the core idea. The difference is scale: modern LLMs have hundreds of billions of parameters, and training requires computing gradients over terabytes of data using thousands of GPUs in parallel.

Deep Learning, CNNs, and RNNs

Deep learning simply means using neural networks with many layers (tens to hundreds). Depth allows the network to learn hierarchical representations—low-level features in early layers, high-level abstractions in later layers. This is analogous to a multi-stage signal processing pipeline.

Convolutional Neural Networks (CNNs) apply learnable spatial filters (convolution kernels) to input data. In VLSI, CNNs have been successfully applied to lithography hotspot detection, where the input is a 2D layout image and the output is a probability map of potential yield issues. They are also used in physical design for routability prediction and IR drop estimation.

Recurrent Neural Networks (RNNs) process sequential data by maintaining a hidden state that acts as memory—similar to how a flip-flop in a sequential circuit stores the current state. Each time step, the network reads a new input and updates its hidden state. RNNs and their improved variant, LSTMs (Long Short-Term Memory), were the foundation of early sequence-to-sequence models for language translation and code generation before transformers superseded them.

Transformers and the Attention Mechanism

The transformer architecture, introduced in the 2017 paper "Attention Is All You Need," is the foundation of every modern large language model (GPT-4, Claude, Gemini, LLaMA). Understanding transformers is essential for evaluating LLM-based RTL generation tools.

The key innovation is the **self-attention mechanism**. Given a sequence of input tokens (words, code tokens, or any discrete symbols), self-attention computes a weighted combination of all tokens in the sequence, where the weights are learned

functions of the token content. This allows the model to capture long-range dependencies without the sequential bottleneck of RNNs.

In EE terms, think of attention as an adaptive, content-dependent interconnect matrix. In a fixed circuit, signal routing is static. In a transformer, the "routing" between tokens is dynamically computed at every layer based on the data itself. This is why transformers excel at language: the meaning of a word depends on context that may be hundreds of tokens away.

Concept	Transformer Term	VLSI Analogy
Input representation	Token embedding	Encoding a signal into a bus representation
Adaptive routing	Self-attention weights	Crossbar switch with learned select lines
Feature extraction	Feed-forward layers	Multi-stage combinational logic
Sequence processing	Positional encoding	Timestamp or address tagging
Output generation	Autoregressive decoding	Shift register outputting one token per cycle

Reinforcement Learning

Reinforcement learning (RL) is the branch of ML concerned with training an **agent** to take **actions** in an **environment** to maximize cumulative **reward**. This maps naturally to many EDA optimization problems:

- **Agent:** The optimization algorithm (e.g., the placement engine)
- **State:** The current design configuration (e.g., macro positions on the floorplan)
- **Action:** A design decision (e.g., place the next macro at coordinates (x, y))
- **Reward:** A quality metric (e.g., negative wirelength, or a composite of timing, area, and congestion)

Google's AlphaChip trained an RL agent to place chip macros by treating the floor-planning problem as a sequential decision game. Synopsys DSO.ai uses RL to navi-

gate the enormous parameter space of synthesis and place-and-route tool options. In verification, RL agents can learn to generate stimulus sequences that maximize coverage metrics, effectively replacing hand-tuned coverage-directed test generation.

Key Terminology Glossary

ML Term	Definition	VLSI Relevance
Inference	Running a trained model on new inputs	Deploying an ML model in your EDA tool flow
Fine-tuning	Adapting a pre-trained model to a specific domain	Training on your company's proprietary RTL
RAG (Retrieval-Augmented Generation)	Augmenting LLM output with retrieved documents	Feeding design specs or coding guidelines into an LLM prompt
Hallucination	Model generating plausible but incorrect output	Syntactically valid but functionally wrong RTL
Epoch	One complete pass through the training data	Analogous to one iteration of a global optimization sweep
Overfitting	Model memorizes training data, fails on new inputs	Like a testbench that only works for specific test vectors
Tokens	Sub-word units that LLMs process	<code>always @(posedge clk)</code> is ~5-7 tokens

Training, Fine-Tuning, and Deployment

Pre-training builds a foundation model on massive, general-purpose data (e.g., trillions of tokens of text and code from the internet). This is enormously expensive—

tens of millions of dollars in compute—and is done by large labs (OpenAI, Google, Meta, Anthropic).

Fine-tuning adapts a pre-trained model to a specific domain or task using a smaller, specialized dataset. NVIDIA's ChipNeMo demonstrated that fine-tuning LLaMA-2 on 24 billion tokens of internal chip design data (RTL, design documents, bug reports) produced a model that significantly outperformed the base model on EDA-specific tasks like code generation and bug summarization. Fine-tuning typically costs 10-100× less than pre-training and can be done on a single 8-GPU server.

Deployment means running inference in a production workflow. For interactive use cases (code assistants, chatbots), latency matters: you need responses in seconds. For batch use cases (analyzing an entire RTL codebase for bugs), throughput matters more. Quantization (reducing model weights from 16-bit to 8-bit or 4-bit) and distillation (training a smaller model to mimic a larger one) are common techniques for making deployment practical without prohibitive GPU costs.

CHAPTER 3

LLMs for RTL Generation

Large language models have emerged as the most immediately accessible AI tool for working RTL designers. Unlike physical design AI (which requires deep integration with P&R tools) or ML-based verification (which requires extensive training infrastructure), LLM-based code generation can be used today with nothing more than a browser or an IDE plugin. This chapter covers how these models work, what is available, and—critically—how to use them effectively and safely.

How LLMs Generate Code

An LLM generates text (including code) by predicting the next token in a sequence, one token at a time. Given the prompt "Write a Verilog module for a 4-bit counter," the model computes a probability distribution over its vocabulary (~32,000–128,000 tokens) and selects the most likely next token. It then appends that token to the context and repeats. This autoregressive process continues until the model produces a stop token or reaches the context window limit.

The **context window** is the maximum number of tokens the model can process at once. Current frontier models advertise context windows from 128K to 1M+ tokens, though effective utilization degrades on long contexts; 200K tokens is a practical ceiling for most production use today. For RTL work, this means you can paste an entire module (or several related modules) into the prompt as context. Larger context windows are directly beneficial for hardware design, where understanding cross-module interfaces and signal dependencies is critical.

KEY CONCEPT: TOKEN PREDICTION, NOT UNDERSTANDING

LLMs do not "understand" Verilog in the way a synthesis tool does. They predict statistically likely token sequences based on patterns learned from training data. This means they can produce syntactically perfect code that is functionally incorrect. Every LLM-generated module must be verified—there are no exceptions to this rule.

General-Purpose Models for RTL

The major general-purpose LLM families from OpenAI, Anthropic, and Google all have significant Verilog and SystemVerilog capability, learned from open-source HDL in their training corpora (primarily GitHub). GitHub Copilot provides inline RTL completion in VS Code and other editors, with model options that evolve over time [23].

These models work best for common, well-documented patterns: standard interfaces (AXI, Wishbone), textbook building blocks (FIFOs, arbiters, UARTs), and straightforward combinational and sequential logic. They struggle with proprietary protocols, non-standard coding styles, and complex microarchitectural logic that is rare in open-source HDL.

Domain-Specific Models

Several research and commercial efforts have produced models specifically optimized for hardware design:

- **ChipNeMo (NVIDIA, 2023):** LLaMA-2 70B fine-tuned on 24 billion tokens of NVIDIA's internal data (RTL, design docs, EDA scripts, bug reports). Achieved 2–5× improvement on chip design tasks versus the base model. Demonstrated the value of domain-adaptive pre-training followed by supervised fine-tuning.
- **RTLCoder (2024):** An open-source model focused on RTL generation, trained on curated Verilog datasets with automated quality filtering. Showed competitive performance with GPT-3.5 at a fraction of the model size.
- **VeriGen (Thakur et al., 2023):** CodeGen-16B fine-tuned on Verilog from GitHub and textbooks. Pioneered evaluation methodologies for RTL-generating LLMs using functional correctness metrics rather than just syntactic validity.
- **ChipChat (2023):** Explored conversational, iterative RTL generation where the model produces code, receives feedback from EDA tools (compilation errors, simulation failures), and iteratively refines its output.

2025-2026 update: The field has shifted from "single model produces RTL" toward *evaluation-driven and agent-assisted flows*. New benchmark suites (for example, CVDP and expanded open RTL benchmarks) indicate that pass@1 functional correctness remains far from solved on realistic tasks, even for strong models [24][25]. Recent work focuses on improving dataset quality (functionally validated training pairs), tool-in-the-loop generation, and retrieval/agent orchestration rather than relying only on larger base models [26][27].

Prompt Engineering for RTL

The quality of LLM-generated RTL depends heavily on prompt quality. Effective prompts for hardware design share several characteristics:

1. **Specify the interface explicitly.** List all ports, widths, directions, and clock/reset conventions. Do not assume the model will infer your interface.
2. **Describe the behavior precisely.** Use unambiguous language. Specify edge cases, reset values, and timing relationships.
3. **Provide coding style constraints.** Specify synchronous vs. asynchronous reset, blocking vs. non-blocking assignments, naming conventions.
4. **Include examples or context.** Paste related modules, interface definitions, or design specifications into the prompt as context.
5. **Request specific structure.** Ask for synthesizable RTL, specify whether you want a Moore or Mealy FSM, indicate whether output should be registered.

Example: Generating an AXI-Lite Slave Register Bank

Production Practice: Treat this pattern as first-draft acceleration. It is useful in day-to-day RTL work, but requires normal lint, CDC, simulation, and review gates before merge.

Here is a concrete example of effective prompting. The prompt specifies the interface, behavior, coding style, and constraints:

```
// Prompt to LLM:
// Write a synthesizable SystemVerilog module for an AXI4-Lite slave
// register bank with the following specifications:
//   - 4 read/write registers at offsets 0x00, 0x04, 0x08, 0x0C
//   - 32-bit data width, byte-addressable
//   - Synchronous active-low reset
//   - Single-cycle read/write response
//   - Use non-blocking assignments for all sequential logic
//   - Follow the naming convention: i_ for inputs, o_ for outputs
```

A well-tuned model will produce something like:

```

module axi_lite_regs (
    input logic      i_aclk,
    input logic      i_aresetn,
    // Write address channel
    input logic [31:0] i_awaddr,
    input logic      i_awvalid,
    output logic      o_awready,
    // Write data channel
    input logic [31:0] i_wdata,
    input logic [3:0]  i_wstrb,
    input logic      i_wvalid,
    output logic      o_wready,
    // Write response channel
    output logic [1:0] o_bresp,
    output logic      o_bvalid,
    input logic      i_bready,
    // Read address channel
    input logic [31:0] i_araddr,
    input logic      i_arvalid,
    output logic      o_arready,
    // Read data channel
    output logic [31:0] o_rdata,
    output logic [1:0] o_rresp,
    output logic      o_rvalid,
    input logic      i_rready
);

logic [31:0] reg_bank [0:3];
logic [1:0]  wr_addr, rd_addr;

assign wr_addr  = i_awaddr[3:2];
assign rd_addr  = i_araddr[3:2];
assign o_bresp  = 2'b00;
assign o_rresp  = 2'b00;

// Write logic
always_ff @(posedge i_aclk) begin
    if (!i_aresetn) begin
        o_awready <= 1'b0;
        o_wready  <= 1'b0;
        o_bvalid  <= 1'b0;
        for (int i = 0; i < 4; i++)
            reg_bank[i] <= 32'h0;
    end
end

```



```

        end else begin
            o_awready <= 1'b1;
            o_wready  <= 1'b1;
            if (i_awvalid && i_wvalid) begin
                for (int b = 0; b < 4; b++)
                    if (i_wstrb[b])
                        reg_bank[wr_addr][b*8 +: 8] <= i_wdata[b*8 +: 8];
                o_bvalid <= 1'b1;
            end else
                o_bvalid <= o_bvalid && !i_bready ? 1'b1 : 1'b0;
        end
    end

    // Read logic
    always_ff @(posedge i_aclk) begin
        if (!i_aresetn) begin
            o_arready <= 1'b0;
            o_rvalid  <= 1'b0;
            o_rdata   <= 32'h0;
        end else begin
            o_arready <= 1'b1;
            if (i_arvalid) begin
                o_rdata <= reg_bank[rd_addr];
                o_rvalid <= 1'b1;
            end else
                o_rvalid <= o_rvalid && !i_rready ? 1'b1 : 1'b0;
        end
    end

end

endmodule

```

This output is plausible and syntactically correct. But notice: it simplifies the AXI4-Lite handshake (awready/wready are held high, and the write address and data channels are assumed to be valid simultaneously). A production implementation would need separate FSMs for each channel. This illustrates a key point: LLM-generated RTL is a *starting point*, not a finished product.

Fine-Tuning on Proprietary HDL

General-purpose models have learned Verilog primarily from open-source repositories, which represent a tiny fraction of the world's RTL. Most high-quality, high-com-

plexity RTL exists behind corporate firewalls. Fine-tuning on your organization's codebase can dramatically improve results for your specific use cases.

NVIDIA's ChipNeMo project provides a practical blueprint:

1. **Data collection:** Aggregate RTL source files, design specifications, verification plans, EDA tool scripts, and bug tracker data. ChipNeMo used 24 billion tokens from NVIDIA's internal repositories.
2. **Domain-adaptive pre-training (DAPT):** Continue pre-training the base model on the domain corpus. This teaches the model your vocabulary, coding conventions, and design patterns.
3. **Supervised fine-tuning (SFT):** Further train on curated prompt-response pairs for specific tasks (e.g., "generate a module from this spec," "explain this code," "find bugs in this block").
4. **Evaluation:** Measure on held-out test sets with functional correctness metrics, not just syntax or BLEU scores.

Limitations You Must Understand

Working with LLMs for RTL requires a clear-eyed understanding of their failure modes:

- **Hallucinations:** Models confidently generate nonexistent SystemVerilog constructs, incorrect port widths, or logic that does not implement the specified behavior. This is not a bug that will be "fixed"—it is an inherent property of statistical token prediction.
- **No timing awareness:** LLMs have no concept of propagation delay, setup/hold times, or clock frequency. They cannot ensure that generated logic meets timing constraints.
- **Verification gap:** The model cannot verify its own output. It does not run simulation, formal checking, or synthesis. You must close this gap with your existing verification infrastructure.
- **Non-determinism:** The same prompt can produce different outputs on different runs. This complicates reproducibility in a design flow that demands it.
- **IP and confidentiality:** Sending proprietary RTL to cloud-based LLMs raises IP concerns. On-premises deployment or API agreements with data protection clauses are essential for production use.

Practical Workflow Integration

The most effective approach today treats LLM-generated RTL as a **first draft** within an established design methodology:

1. **Specification:** Write a clear natural-language or structured specification for the module.
2. **Generation:** Use an LLM (via IDE plugin or API) to produce an initial implementation.
3. **Review:** The designer reviews the generated code, correcting errors and refining the implementation. This step is non-negotiable.
4. **Lint and compile:** Run the code through your standard linting and synthesis flow. Feed errors back to the LLM for correction (the ChipChat iterative approach).
5. **Verification:** Run the module through your full verification flow—simulation, formal, and code coverage—exactly as you would for hand-written RTL.
6. **Commit:** Only code that passes all quality gates enters the design repository.

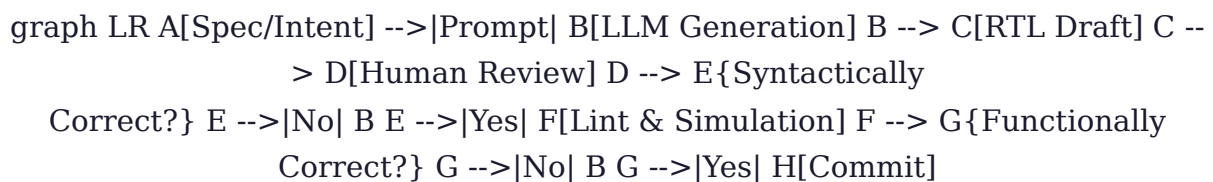


Figure 3.1: The AI-Assisted RTL Design Loop

The productivity gain from LLMs is not in eliminating the verification step—it is in accelerating the time from specification to first compilable draft. For well-specified blocks, this can reduce initial coding time from days to minutes.

LLM RTL Acceptance Checklist

- **Interface correctness:** Port list, widths, directions, reset polarity, and clocking scheme match the spec.
- **Style conformance:** Blocking/non-blocking use, reset style, naming conventions, and synthesis subset rules are enforced.

- **Static quality gates:** Lint, CDC/RDC checks, and compile/elaboration pass with no unreviewed waivers.
- **Behavioral confidence:** Directed tests plus constrained-random regressions show expected behavior and no corner-case regressions.
- **Sign-off readiness:** Coverage thresholds, assertions, and code review comments are resolved before commit.

CHAPTER 4

AI-Assisted RTL Code Review and Bug Detection

Every experienced RTL designer knows that certain classes of bugs are maddeningly difficult to catch through simulation alone. Clock domain crossings (CDCs) that fail once per billion cycles, unintended latches inferred from incomplete case statements, and finite state machines with unreachable states—these are the defects that escape verification and surface in silicon. Traditional lint tools catch the obvious cases through pattern matching on predefined rules. ML-augmented analysis can go further, learning to recognize subtle patterns that no human has written a rule for.

Traditional Linting: Strengths and Limits

Tools like Synopsys SpyGlass, Siemens Questa AutoCheck, and Cadence JasperGold CDC operate on rule-based pattern matching. They parse the RTL into an abstract syntax tree (AST) or a structural netlist, then apply hundreds of predefined checks: combinational loops, multiple drivers, uninitialized registers, sensitivity list mismatches, and CDC violations.

These tools are mature and invaluable, but they share fundamental limitations:

- **Fixed rule sets:** They can only detect what they have been programmed to detect. Novel bug patterns, design-specific anti-patterns, or subtle cross-module interactions that no rule anticipates will pass through.
- **High false-positive rates:** Conservative rules produce large volumes of warnings. Industry surveys and user reports frequently show substantial false-positive overhead, leading to "lint fatigue" where real issues are missed in the noise [9].
- **No learning from history:** Traditional tools do not improve from exposure to your design team's historical bugs. The same types of errors recur project after project, but the linter does not adapt.

KEY CONCEPT: FROM RULES TO LEARNED PATTERNS

ML-augmented code analysis shifts from "match this predefined pattern" to "this code is statistically anomalous relative to correct code in the training corpus." This allows detection of bug classes that no human engineer has explicitly codified as rules, including design-specific anti-patterns unique to your codebase.

RTL Bug Patterns That ML Can Detect

Clock Domain Crossings (CDCs)

CDC bugs are among the most dangerous defects in digital design because they are inherently non-deterministic—metastability-induced failures depend on the exact phase relationship between asynchronous clocks at the moment of a transition. Traditional CDC tools verify structural synchronization (e.g., presence of a 2-FF synchronizer), but ML models can learn higher-level patterns: whether the synchronization scheme is appropriate for the data type (single-bit vs. multi-bit), whether gray coding is used where required, and whether the associated handshake protocol is complete.

Academic work has shown that graph neural networks (GNNs) trained on annotated circuit/netlist datasets can improve bug-pattern detection quality over rule-only baselines on curated benchmarks, while reducing noise from recurring structural patterns [22].

Unintended Latches

An incomplete `if-else` or `case` statement in combinational logic causes synthesis to infer a latch—almost always a bug. While basic lint rules catch the simplest cases, ML models can detect more subtle variants: latches inferred through complex conditional nesting, latches hidden behind generate blocks, or combinational paths that are *technically* complete but functionally equivalent to latched behavior due to constant propagation.

FSM Deadlocks and Unreachable States

A finite state machine with a state that has no exit transition (deadlock) or no entry path (unreachable) is a common design error, especially in complex protocol imple-

mentations. Traditional tools flag obvious cases, but ML approaches that model the FSM as a graph and apply GNN-based reachability analysis can detect deadlocks that depend on specific input sequences or multi-cycle conditions that are difficult to express as static lint rules.

Reset Domain Errors

In multi-reset designs, registers that are reset by different signals can create ordering-dependent behavior. ML models trained on correct reset structures can flag anomalous reset topologies—for example, a control register reset by one domain driving datapath logic reset by another, without proper synchronization between the reset release sequences.

ML Approaches for RTL Analysis

Graph Neural Networks on Circuit Netlists

RTL and netlists are inherently graph-structured: modules are nodes, connections are edges, and hierarchical composition creates nested subgraphs. This makes graph neural networks a natural fit for structural analysis.

A GNN-based bug detection system typically works as follows:

1. **Parse** the RTL into an elaborated netlist or an intermediate representation such as FIRRTL or Yosys's RTLIL.
2. **Construct a graph** where nodes represent cells (gates, flip-flops, muxes) and edges represent connections. Node features encode cell type, bitwidth, and other attributes.
3. **Apply message-passing layers** that propagate information along edges, allowing each node to accumulate context from its neighborhood (analogous to signal propagation through the circuit).
4. **Classify** nodes, edges, or subgraphs as "normal" or "potentially buggy" using a trained classifier head.

```
graph TD
  A[RTL Source] -->|Parsing| B[Netlist Graph]
  B --> C[Feature Extraction]
  C --> D[GNN Model Inference]
  D --> E["{Anomaly Score}"]
  E -->|High| F[Flag Potential Bug]
  E -->|Low| G[Pass]
  H[Historical Bug Data] -->|Training| D
```

Figure 4.1: GNN-Based Bug Detection Pipeline

```

# Simplified GNN-based RTL anomaly detection pipeline
import torch
from torch_geometric.nn import GCNConv
from torch_geometric.data import Data

class RTLBugDetector(torch.nn.Module):
    def __init__(self, in_features, hidden_dim, num_classes):
        super().__init__()
        self.conv1 = GCNConv(in_features, hidden_dim)
        self.conv2 = GCNConv(hidden_dim, hidden_dim)
        self.classifier = torch.nn.Linear(hidden_dim, num_classes)

    def forward(self, data):
        x, edge_index = data.x, data.edge_index
        x = self.conv1(x, edge_index).relu()
        x = self.conv2(x, edge_index).relu()
        return self.classifier(x)

# Node features encode: cell type, bitwidth, clock domain,
# reset domain, combinational depth from primary input

```

Anomaly Detection with Autoencoders

An alternative approach trains an autoencoder to reconstruct features of "correct" RTL code (extracted from a vetted, production-quality codebase). During inference, code regions with high reconstruction error are flagged as anomalous. This unsupervised approach has the advantage of not requiring labeled bug data—it learns what "normal" looks like and flags deviations.

LLM-Based Code Review

General-purpose and fine-tuned LLMs can perform natural-language code review on RTL, flagging potential issues and explaining their reasoning. While less precise than structural analysis, this approach can catch higher-level design issues: protocol violations, inconsistencies between code and comments, and deviations from design specifications when the spec is provided as context.

Tools and Research

Tool / Project	Approach	Status
Synopsys SpyGlass + AI	ML-enhanced lint rules with false-positive suppression	Commercial (production)
Cadence Verisium AI	ML-based regression failure triage and root-cause analysis	Commercial (production)
CircuitNet (PKU)	Open-source dataset and GNN benchmarks for circuit analysis	Academic / Open-source
DeepGate (HKUST)	GNN-based gate-level logic reasoning and prediction	Academic
RTLLM (Various)	LLM-based RTL generation benchmarks with correctness evaluation	Academic

Case Study: ML-Augmented CDC Verification

Consider an illustrative scenario: an SoC integration team managing 15 clock domains with over 800 CDC crossings. The traditional CDC flow produces 2,400 warnings, of which approximately 1,500 are false positives that engineers must manually triage—a process that takes two engineers three weeks per project milestone.

By training a gradient-boosted classifier on historical triage data (five projects, ~12,000 labeled warnings), the team built a model that predicts false positives with high precision and recall. Integrating this as a post-filter reduces actionable warnings from 2,400 to approximately 600, cutting triage time from three weeks to four days. This style of workflow is common in practical ML triage deployments [11][12].

The training pipeline is straightforward:

```
# Feature extraction from SpyGlass CDC warnings
features = [
    'crossing_type',      # single-bit, multi-bit, bus
    'sync_scheme',       # 2FF, handshake, FIFO, none
    'src_clock_freq',    # source clock frequency
    'dst_clock_freq',    # destination clock frequency
    'freq_ratio',        # ratio of src/dst frequencies
    'fanout_count',      # number of destination registers
    'reconvergence_depth', # depth of reconvergent paths
    'module_hierarchy',  # encoded hierarchical path
    'signal_name_pattern', # TF-IDF of signal naming conventions
    'historical_waiver',  # was a similar warning waived before?
]

# Labels: 0 = true bug, 1 = false positive (from historical triage)
from sklearn.ensemble import GradientBoostingClassifier
model = GradientBoostingClassifier(n_estimators=200, max_depth=6)
model.fit(X_train, y_train)
```

Integration into CI/CD for Hardware

Modern RTL development teams increasingly adopt continuous integration practices borrowed from software engineering. ML-based code analysis fits naturally into this workflow:

1. **Pre-commit hooks:** Run lightweight LLM-based code review on changed files before they enter the repository. Flag potential issues in the merge request comments.
2. **Nightly regression:** Run full structural analysis (lint + ML post-filtering) on the entire design. Track warning trends over time to catch regressions early.
3. **Milestone gates:** Require ML-augmented CDC and lint reports as sign-off criteria at design milestones (RTL freeze, synthesis handoff, tapeout).
4. **Feedback loop:** When engineers triage warnings and file waivers, feed this data back into the ML model as new training labels. The model improves with each project cycle.

The goal is not to replace expert review—it is to focus expert attention where it matters most. An ML filter that eliminates 60% of false positives frees your senior engineers to spend their time on the warnings that are most likely to be real bugs.

In the next chapter, we will move from static code analysis to dynamic verification—how AI and ML are transforming functional simulation, coverage closure, and test-bench generation.

CHAPTER 5

AI for Functional Verification

The Verification Bottleneck

Functional verification often dominates project effort on modern SoCs [9]. This is not a tooling failure—it is a direct consequence of exponential design complexity. A block with n state bits has 2^n reachable states; a system integrating hundreds of such blocks, each with asynchronous interfaces and configurable modes, creates a state space that no amount of directed testing can exhaust. The industry's standard response—constrained-random verification within a UVM testbench—has been remarkably successful, but it is reaching its limits. Coverage closure on complex designs now takes months, and the cost of a single missed corner case is significant.

AI and ML offer a fundamentally different approach: instead of relying on hand-crafted stimulus constraints and coverage-driven heuristics, learned policies can explore the state space more efficiently, close coverage gaps faster, and even generate assertions and testbench components automatically.

Coverage-Driven Verification and Its Challenges

In a standard UVM flow, the verification engineer defines **functional coverage groups**—cross products of signal values, FSM state transitions, and protocol events—that represent the behaviors the design must exhibit. The constrained-random testbench generates stimulus that satisfies specified constraints, and the coverage database tracks which coverage points have been hit. When coverage stalls, the engineer manually writes directed tests, tweaks constraint weights, or adds new coverage-closing sequences.

This process has several pain points that ML can address:

- **Constraint engineering is manual and iterative.** Writing constraints that are both legal (satisfying protocol rules) and useful (reaching uncovered states) requires deep design knowledge and extensive trial-and-error.

- **Coverage closure follows a long tail.** The first 80% of coverage points are typically hit quickly. The remaining 20% can take 5–10× longer because they require rare stimulus combinations.
- **Feedback is slow.** Each simulation run takes minutes to hours. The engineer analyzes the coverage report, modifies constraints, and reruns—a cycle measured in days.
- **Cross-coverage explosion.** Crossing k coverage dimensions with m values each produces m^k bins. Most are unreachable, but the engineer must manually exclude them or wait for tools to prune.

KEY CONCEPT: VERIFICATION AS EXPLORATION

Coverage closure is fundamentally a search problem: navigate a vast state space to visit specific target states. This maps directly to reinforcement learning, where an agent learns a policy to maximize a reward signal (coverage hits) through sequential decisions (stimulus choices). RL agents can learn to reach hard-to-cover states orders of magnitude faster than uniform random exploration.

Reinforcement Learning for Coverage Closure

The application of RL to coverage-driven verification works as follows. The verification environment is treated as the RL **environment**. The RL **agent** controls the stimulus generation—either by selecting constraint parameters, choosing from a menu of pre-defined test scenarios, or directly generating input sequences. The **reward** is defined in terms of coverage: new coverage points hit, or the incremental change in coverage metrics after each simulation episode.

A typical architecture uses Proximal Policy Optimization (PPO) or Soft Actor-Critic (SAC) as the RL algorithm, with the state representation encoding the current coverage database status, recent stimulus history, and design configuration registers. The agent learns which stimulus distributions are most likely to reach uncovered states, effectively replacing weeks of manual constraint tuning with an automated policy that improves with each simulation run.

Practical Example: RL for Coverage-Driven Random Verification

Consider a PCIe transaction-layer verification environment with 450 functional coverage points across transaction types, sizes, ordering rules, and error injection scenarios. After two weeks of standard constrained-random simulation, coverage has stalled at 87%. The remaining 13%—59 coverage points—involves rare combinations of out-of-order completions, specific TLP size/address alignments, and error-recovery sequences.

```

# RL agent for PCIe coverage closure
import gymnasium as gym
import numpy as np
from stable_baselines3 import PPO

class PCIECoverageEnv(gym.Env):
    def __init__(self, sim_interface, coverage_db):
        super().__init__()
        self.sim = sim_interface
        self.cov = coverage_db

        # Action: select transaction type, size, address pattern,
        # ordering mode, error injection flag
        self.action_space = gym.spaces.MultiDiscrete([8, 6, 4, 3, 2])

        # State: coverage bitmap + recent transaction history
        self.observation_space = gym.spaces.Box(
            low=0, high=1,
            shape=(450 + 64,), dtype=np.float32
        )

    def step(self, action):
        txn = self._decode_action(action)
        self.sim.drive_transaction(txn)
        self.sim.run_cycles(1000)

        new_cov = self.cov.get_bitmap()
        new_hits = np.sum(new_cov & ~self.prev_cov)
        self.prev_cov = new_cov

        reward = new_hits * 10.0
        total_coverage = np.mean(new_cov)
        done = total_coverage >= 0.99

        obs = np.concatenate([new_cov, self.history])
        return obs, reward, done, False, {}

    def _decode_action(self, action):
        # Map discrete action indices to PCIe transaction parameters
        return {
            'type': ['MRd', 'MWr', 'IO Rd', 'IO Wr',
                    'Cfg Rd', 'Cfg Wr', 'Msg', 'Cpl'][action[0]],
            'size': [1, 2, 4, 16, 64, 256][action[1]],

```

```

        'addr_pattern': ['aligned', 'unaligned',
                        '4k_boundary', 'wrap'][action[2]],
        'ordering': ['strong', 'relaxed', 'no_snoop'][action[3]],
        'inject_error': bool(action[4]),
    }

# Train the agent
env = PCIeCoverageEnv(sim_interface, coverage_db)
agent = PPO("MlpPolicy", env, verbose=1, learning_rate=3e-4)
agent.learn(total_timesteps=50000)

```

In published case studies using this approach, RL-guided stimulus often reaches hard-to-cover bins faster than undirected constrained-random methods, with the largest gains on long-tail coverage points. The key insight is that the agent learns to *focus* stimulus on unexplored regions of the state space rather than uniformly sampling the entire distribution [22].

Assertion Mining from Simulation Traces

Writing SystemVerilog Assertions (SVAs) is tedious and error-prone. A verification engineer must translate informal protocol knowledge into precise temporal logic—and the number of assertions needed grows with design complexity. ML offers an alternative: automatically mine likely invariants and temporal properties from simulation traces.

The process works by observing signal behavior across thousands of simulation cycles and identifying patterns that hold consistently:

- **Invariant detection:** Identify relationships like `grant |-> ##[1:3] ack` that hold across all observed traces.
- **Temporal pattern mining:** Discover sequences such as "request is always followed by grant within 5 cycles when arbiter is not in reset."
- **Anomaly-based assertions:** Train a model on "golden" simulation runs, then flag future runs that deviate from learned temporal patterns as potential bugs.

Tools like Synopsys VC Formal's assertion synthesis and Cadence JasperGold's coverage unreachable analysis already incorporate variants of this approach. The mined assertions serve as a starting point; the verification engineer reviews, refines, and promotes them into the formal testbench.

Intelligent Testbench Generation with LLMs

LLMs can accelerate testbench development by generating UVM component scaffolding, cocotb test classes, and constrained-random sequences from natural-language specifications. This is particularly effective for standard protocol verification environments where the structure is well-established and extensively documented in the LLM's training data.

```

# LLM-generated cocotb test for an AXI4 FIFO
import cocotb
from cocotb.triggers import RisingEdge, ClockCycles
from cocotb.clock import Clock
from cocotbext.axi import AxiStreamSource, AxiStreamSink, AxiStreamBus

@cocotb.test()
async def test_fifo_backpressure(dut):
    """Verify FIFO handles sink backpressure without data loss."""
    clock = Clock(dut.clk, 10, units="ns")
    cocotb.start_soon(clock.start())

    source = AxiStreamSource(
        AxiStreamBus.from_prefix(dut, "s_axis"), dut.clk, dut.rst
    )
    sink = AxiStreamSink(
        AxiStreamBus.from_prefix(dut, "m_axis"), dut.clk, dut.rst
    )

    dut.rst.value = 1
    await ClockCycles(dut.clk, 10)
    dut.rst.value = 0
    await ClockCycles(dut.clk, 5)

    sink.set_pause_generator(
        iter([1, 1, 1, 0, 0, 1, 0, 1, 1, 0] * 100)
    )

    test_data = [bytes([i % 256] * 64) for i in range(200)]
    for frame in test_data:
        await source.send(frame)

    await ClockCycles(dut.clk, 5000)

    received = []
    while not sink.empty():
        received.append(bytes(await sink.recv()))

    assert len(received) == len(test_data), \
        f>Data loss: sent {len(test_data)}, received {len(received)}"
    for i, (sent, got) in enumerate(zip(test_data, received)):
        assert sent == got, f"Mismatch at frame {i}"

```

The LLM handles the boilerplate—clock setup, reset sequencing, bus instantiation—while the engineer specifies the scenario (backpressure patterns, data volumes, assertion criteria). For a complex UVM environment, the time savings on sequence library scaffolding alone can be substantial.

Commercial Tool Landscape

Tool	Vendor	ML Capability
VCS with Verdi ML	Synopsys	ML-driven regression optimization, smart exclusion of unreachable coverage, test ranking by predicted coverage contribution
Verisium AI	Cadence	Automated failure triage, root-cause clustering, regression suite optimization, coverage convergence prediction
Questa inFact	Siemens EDA	Graph-based intelligent testbench automation, portable stimulus with ML-guided constraint solving
VC Formal AI	Synopsys	ML-guided proof strategies, assertion synthesis, coverage unreachability analysis

Commercial verification platforms now include ML-guided regression planning, failure clustering, and coverage analysis. Synopsys emphasizes AI-assisted verification space optimization and formal coverage analysis, while Cadence Verisium focuses on AI-driven triage/debug workflows [11][12][15].

The highest-ROI application of ML in verification today is not generating stimulus—it is triaging the results. On large SoC projects with thousands of nightly regression tests, ML-based failure clustering and root-cause prediction can save weeks of engineering time per milestone.

CHAPTER 6

AI for Formal Verification

Formal Methods Overview

Formal verification proves (or disproves) that a design satisfies a set of properties for *all possible* input sequences, not just the ones that happen to appear in a simulation testbench. This exhaustive guarantee makes formal methods the gold standard for safety-critical logic: bus protocol compliance, arbiter fairness, FIFO overflow prevention, and security-sensitive control paths.

The three primary formal techniques in production use are:

- **Model checking:** Exhaustively explores the reachable state space of a design to verify that a temporal property (expressed in SVA, PSL, or CTL/LTL) holds on every reachable state. Tools include Synopsys VC Formal, Cadence JasperGold, and Siemens Questa Formal.
- **Equivalence checking:** Proves that two representations of a design (e.g., RTL vs. gate-level netlist, or pre- and post-ECO RTL) are functionally identical. This is the most widely deployed formal technique, used in every tapeout flow.
- **Property checking (assertion-based):** Verifies specific SVA/PSL properties against the RTL. Unlike full model checking, property checking can be applied incrementally to individual properties, making it more scalable for targeted verification.

The State Space Explosion Problem

The fundamental challenge of model checking is state space explosion. A design with n state-holding elements has up to 2^n reachable states. A modest 32-bit register file with 16 entries contains 512 state bits, yielding a theoretical state space of 2^{512} —vastly more than the number of atoms in the observable universe. Symbolic model checking (using Binary Decision Diagrams or SAT/SMT solvers) and bounded model checking (BMC) have pushed practical limits significantly, but even these techniques routinely fail to converge on complex control logic with deep sequential depth or wide datapaths.

This is exactly where ML can help: not by replacing the mathematical machinery of formal provers, but by guiding their search strategies, decomposing hard problems into tractable subproblems, and learning abstractions that reduce the effective state space.

KEY CONCEPT: ML AS A GUIDE, NOT A REPLACEMENT

ML does not replace the formal proof engine—the proof itself must still be mathematically rigorous. Instead, ML acts as a heuristic guide that helps the prover explore the state space more intelligently. Think of it as the difference between brute-force search and A* search: the guarantee of correctness is the same, but the path to finding the proof (or counterexample) is dramatically shorter.

ML-Guided Search Strategies

Learned Abstractions

Abstraction is the primary technique for taming state space explosion: replace parts of the design with simpler models that preserve the properties of interest while reducing complexity. Traditionally, engineers manually identify which signals to abstract and how. ML can automate this by learning which components of a design are irrelevant to a given property.

A GNN trained on the circuit graph can predict, for each node, the probability that it influences the property under verification. Nodes with low influence scores are candidates for abstraction—replaced with unconstrained inputs or simplified models. Published research reports meaningful runtime improvements from learned abstraction strategies, especially on deep sequential properties [22].

Guided SAT/SMT Solving

Modern formal verification engines are built on SAT and SMT solvers. These solvers make branching decisions when exploring the search space, and the order of these decisions dramatically affects performance—a good branching heuristic can mean the difference between solving in seconds and timing out after hours.

ML-based branching heuristics train a model on solved instances to predict which variable assignments are most likely to lead to a solution (or counterexample)

quickly. NeuroSAT-style and follow-on neural-guided SAT approaches have shown practical gains on selected benchmark classes when integrated with conventional solvers [14].

Proof Decomposition

For properties that are too complex to verify monolithically, decomposition splits the proof into smaller subproblems that can be solved independently. An ML model can learn decomposition strategies from successfully proven properties: which assume-guarantee patterns work, where to insert cut points in the design hierarchy, and how to partition temporal properties into bounded segments.

LLMs for Assertion Generation

One of the most immediately practical applications of AI in formal verification is generating SVA properties from natural-language specifications. Verification engineers spend significant time translating English-language protocol descriptions into formal temporal logic—a process that is both tedious and error-prone.

Example: Natural Language to SVA

Research Result + Production Practice: LLM-generated SVA drafts can improve authoring speed, but must be validated by experienced formal engineers before sign-off.

Consider this excerpt from an arbiter specification:

"When a requestor asserts its request signal, it must receive a grant within 8 clock cycles. Once granted, the requestor must hold its request low for at least 2 cycles before re-requesting. No two requestors may be granted simultaneously."

An LLM fine-tuned on SVA corpora can generate:

```

// Property 1: Grant latency bound
property grant_latency(req, gnt);
    @(posedge clk) disable iff (!rst_n)
        $rose(req) |-> ##[1:8] gnt;
endproperty

// Property 2: Minimum re-request interval
property rerequest_interval(req, gnt);
    @(posedge clk) disable iff (!rst_n)
        $rose(gnt) |-> !req[*2] ##1 true;
endproperty

// Property 3: Mutual exclusion
property mutex_grant;
    @(posedge clk) disable iff (!rst_n)
        $onehot0(gnt_bus);
endproperty

// Bind assertions to all requestors
generate for (genvar i = 0; i < NUM_REQ; i++) begin : gen_arb_props
    assert_grant_lat : assert property (
        grant_latency(req[i], gnt[i])
    ) else $error("Grant latency violated for req %0d", i);

    assert_rereq : assert property (
        rerequest_interval(req[i], gnt[i])
    ) else $error("Re-request interval violated for req %0d", i);
end endgenerate

assert_mutex : assert property (mutex_grant)
    else $error("Multiple simultaneous grants detected");

```

This output is close to correct but requires expert review. Common LLM mistakes in SVA generation include incorrect use of `##` delay ranges, confusion between overlapping and non-overlapping implication (`|->` vs. `|=>`), and missing edge cases in the temporal logic. The engineer must verify that the generated properties actually capture the intended specification—but starting from a reasonable draft is far faster than writing from scratch.

SVA Review Checklist for LLM Output

- **Implication semantics:** Confirm use of `| ->` versus `| =>` matches intended cycle alignment.
- **Reset behavior:** Check `disable iff` conditions, reset polarity, and reset release edge cases.
- **Temporal bounds:** Validate delay windows (e.g., `##[1:8]`) against protocol timing requirements.
- **Vacuity resistance:** Ensure antecedents are reachable and properties do not pass trivially.
- **Binding coverage:** Verify generated properties are instantiated for all relevant channels/requestors and reviewed with counterexamples.

ML-Guided Bounded Model Checking

Bounded model checking (BMC) unrolls the design for k time steps and encodes the property check as a SAT problem. The bound k is critical: too small and real bugs are missed; too large and the SAT instance becomes intractable. Traditional approaches increment k linearly.

ML can improve BMC in two ways:

1. **Bound prediction:** A model trained on historical verification data predicts the minimum bound needed to expose a bug for a given property/design pair. This avoids wasting time on shallow bounds when the bug requires deep unrolling, and avoids the exponential cost of unnecessarily deep bounds when the bug is shallow.
2. **Interpolant learning:** Craig interpolants are used in unbounded model checking to generalize from bounded proofs. ML models can predict likely interpolants, accelerating the convergence of k-induction and interpolation-based provers.

Technique	Traditional Approach	ML-Enhanced Approach	Typical Speedup
Abstraction	Manual cone-of-influence	GNN-predicted property relevance	2–10×
SAT branching	VSIDS heuristic	Neural branching predictor	1.5–3×
BMC bound selection	Linear increment	Predicted optimal bound	2–5×
Proof decomposition	Manual assume-guarantee	Learned decomposition	3–8×
Assertion writing	Manual SVA authoring	LLM generation + review	2–4× (engineer time)

Research Frontiers: Neural Theorem Provers

The most ambitious research direction is training neural networks to act as theorem provers for hardware properties. Projects at Google DeepMind, Stanford, and MIT are exploring models that can learn proof strategies end-to-end: given a circuit and a property, directly output a proof (or counterexample) without relying on traditional SAT/SMT machinery.

These are still far from production readiness. Current neural provers work only on small circuits (hundreds of gates) and restricted property classes. But they point toward a future where formal verification scales to entire SoC subsystems—a capability that would fundamentally change how chips are designed and validated.

Practical Considerations and Limitations

Engineers evaluating ML-augmented formal verification should keep several realities in mind:

- **The proof guarantee is unchanged.** ML heuristics guide the search; they do not weaken the formal guarantee. If the prover says a property holds, it holds for all states—the ML component only affects how quickly the answer is found.
- **Training data requirements.** ML-guided formal methods require a corpus of solved verification instances to train on. This means they work best in organizations with large historical databases of formal verification runs.
- **Property quality is still paramount.** No amount of ML acceleration helps if the properties are wrong. LLM-generated SVA must be reviewed by engineers who understand the specification deeply.
- **Integration complexity.** Most ML-guided formal techniques are currently available only in commercial tools or as research prototypes. Integrating custom ML heuristics with production formal engines requires deep tool API knowledge.

Formal verification is where AI in EDA meets its hardest test. The bar is absolute correctness, not statistical improvement. ML's role here is not to relax that bar but to make it achievable on designs that are currently too complex to verify formally.

CHAPTER 7

Tools, Frameworks, and Workflows

Commercial Tools Survey

The three major EDA vendors have all made significant investments in AI/ML capabilities over the past five years. Understanding what each platform offers—and where the marketing claims diverge from engineering reality—is essential for making informed procurement and adoption decisions.

Synopsys DSO.ai and the Synopsys.ai Platform

Synopsys DSO.ai (Design Space Optimization AI) was one of the first production-deployed RL systems in EDA. It operates on the synthesis and place-and-route parameter space—the hundreds of tool options, floorplan constraints, clock tree strategies, and optimization passes that a physical design engineer manually configures across multiple iterations.

DSO.ai uses a reinforcement learning agent that:

1. Explores the tool configuration space by running synthesis and P&R with different parameter combinations.
2. Evaluates results against user-defined objectives (timing, area, power, congestion).
3. Learns which parameter regions are most promising and focuses exploration accordingly.
4. Runs autonomously on compute clusters, exploring tens to hundreds of configurations in parallel.

Synopsys publicly reported in 2023 that DSO.ai had reached over 100 commercial tapeouts and showed up to 25% lower total power in published customer examples; newer counts should be validated against the latest Synopsys disclosures. The broader Synopsys.ai platform extends this to the full design flow [15].

- **VSO.ai:** Verification Space Optimization—applies similar RL-based exploration to verification parameter spaces.
- **TSO.ai:** Test Space Optimization for ATPG and DFT configuration.
- **ASO.ai:** Analog Space Optimization for circuit sizing and layout.

Cadence Cerebrus

Cadence Cerebrus is a full-stack AI engine for digital implementation. It integrates with Genus (synthesis) and Innovus (P&R) to autonomously optimize the implementation flow. Like DSO.ai, it uses reinforcement learning and Bayesian optimization to explore the design space, but Cadence emphasizes the tight coupling between Cerebrus and the underlying tool engines—the AI does not just select parameters but actively guides the internal optimization algorithms.

Key capabilities:

- **Multi-objective optimization:** Simultaneously optimizes timing, power, area, and routability with user-defined weights.
- **Transfer learning:** Knowledge from one design block or technology node can be transferred to accelerate optimization of similar blocks, reducing the number of exploration runs needed.
- **Workload-aware scheduling:** Distributes parallel exploration runs across compute clusters with intelligent resource allocation.

Cadence reports significant schedule and PPA benefits in customer deployments, with final results dependent on design class, constraints, and compute budget [16].

Siemens EDA Solido AI

Solido AI targets analog and mixed-signal design, where the verification challenge is characterizing circuit performance across process, voltage, and temperature (PVT) corners. Traditional Monte Carlo simulation can require very large run counts to achieve statistical confidence on yield estimates. Solido applies ML-based importance sampling and surrogate modeling to reduce required simulation effort in many flows [17].

Applications include:

- **Variation-aware design:** ML surrogate models predict circuit performance across the PVT space, enabling rapid design centering.

- **High-sigma yield analysis:** Importance sampling guided by ML models efficiently estimates failure rates at 6σ and beyond—critical for SRAM bit cells and I/O circuits.
- **Library characterization:** ML accelerates the characterization of standard cell libraries across corners, reducing a traditionally multi-week process to days.

Open-Source and Academic Tools

OpenROAD with ML Plugins

OpenROAD is an open-source RTL-to-GDSII flow developed under DARPA's IDEA program [18]. It includes ML-based components for:

- **Macro placement:** An RL-based placer inspired by AlphaChip, integrated into the OpenROAD flow.
- **CTS optimization:** ML-predicted clock skew for early-stage clock tree evaluation.
- **Congestion prediction:** CNN-based models that predict routing congestion from placement, enabling early floorplan evaluation without running the router.

```
# Running OpenROAD with ML-enhanced macro placement
openroad -exit -no_init <<EOF
read_lef "platforms/sky130/lef/sky130_fd_sc_hd.tlef"
read_def "designs/gcd/gcd.def"

# Enable RL-based macro placement
set_macro_placement_options -method reinforcement_learning \
    -num_episodes 200 \
    -reward_function "wirelength+congestion"
macro_placement

detailed_placement
global_route
detailed_route
write_def "results/gcd_placed_routed.def"
EOF
```

CIRCT and MLIR for Hardware

CIRCT (Circuit IR Compilers and Tools), built on LLVM's MLIR infrastructure, provides a modular compiler framework for hardware design languages [19]. While not an AI tool itself, CIRCT's multi-level IR representation makes it an excellent substrate for ML-based analysis and transformation. Research groups have built GNN-based optimization passes, ML-guided scheduling algorithms, and LLM-driven code transformation pipelines on top of CIRCT's IR.

Google's Chip Design RL

Google's AlphaChip (published in Nature, 2021; open-sourced 2023) demonstrated that RL agents can produce chip macro placements competitive with or superior to human experts [13][20]. The open-source release includes the training infrastructure, environment, and pre-trained models. While most teams will use this as a reference implementation rather than a production tool, it provides an invaluable starting point for teams building custom RL-based optimization.

HDL Generation Frameworks

Framework	Description	AI/ML Relevance
Chisel (Berkeley/SiFive)	Scala-embedded HDL with parameterized generators	LLMs can generate Chisel generators, enabling meta-level design automation
SpinalHDL	Scala-based HDL with strong type safety	Rich type system helps constrain LLM output
Amaranth (nMigen)	Python-based HDL for digital design	Python familiarity makes LLM generation more reliable
PyRTL	Python hardware design library	Simple semantics reduce LLM hallucination risk

Integration Patterns

Embedding AI into an existing EDA flow requires careful architectural decisions. Three integration patterns dominate:

- 1. Pre-processing (AI before EDA):** Use ML models to predict outcomes and guide tool configuration *before* running the traditional flow. Examples: congestion prediction before P&R, coverage prediction before simulation. Low risk, easy to adopt.
- 2. In-loop (AI inside EDA):** ML models operate within the EDA tool's optimization loop, making real-time decisions. Examples: DSO.ai's RL agent, ML-guided SAT branching. Higher impact but requires tool vendor support or deep API access.
- 3. Post-processing (AI after EDA):** Use ML to analyze EDA tool outputs—triage warnings, cluster failures, predict remaining bugs. Examples: Verisium AI regression triage, ML-filtered lint results. Lowest barrier to entry; can be built entirely in-house.

Build vs. Buy Decision Framework

Factor	Favors Build (In-House)	Favors Buy (Commercial)
Problem specificity	Unique to your design flow or methodology	Common across the industry
Data advantage	You have large proprietary datasets	Public or vendor-collected data is sufficient
ML expertise	Team has ML engineering capability	No in-house ML expertise
Integration depth	Needs tight coupling with custom tools	Works within standard vendor tool flows
Time to value	6-18 months acceptable	Need results within 1-3 months
Maintenance	Team can maintain long-term	Prefer vendor-managed updates

For most teams, the practical answer is a hybrid: buy commercial AI-enhanced EDA tools for physical design and verification (where vendor integration is deep and switching costs are high), and build in-house solutions for workflow-specific tasks like lint triage, code review automation, and design-specific coverage prediction where your proprietary data provides a unique advantage.

Data Infrastructure

Every ML application depends on data, and most semiconductor companies are sitting on vast quantities of valuable training data that is poorly organized and underutilized. Building a data infrastructure for AI-driven EDA requires:

- **Design history:** Version-controlled RTL with associated synthesis results, timing reports, and tapeout outcomes. Link commits to bug reports and respin causes.
- **Verification data:** Coverage databases, regression results, failure logs, and manual triage labels. This is the training data for coverage prediction and failure triage models.
- **Tool run metadata:** Synthesis and P&R tool configurations, runtime statistics, and quality-of-result metrics across projects. This trains design space exploration models.
- **Engineer feedback:** Lint waiver decisions, code review comments, and formal property triage results. This is the labeled data that supervised ML models need.

```
# Schema for a minimal verification data warehouse
CREATE TABLE regression_runs (
  run_id      UUID PRIMARY KEY,
  project     VARCHAR(64),
  design_rev  VARCHAR(40),  -- git SHA
  test_name   VARCHAR(256),
  seed        BIGINT,
  status      VARCHAR(16),  -- PASS, FAIL, TIMEOUT
  runtime_sec FLOAT,
  coverage_json JSONB,      -- per-covergroup hit counts
  failure_sig VARCHAR(512), -- failure signature hash
  root_cause  VARCHAR(256), -- engineer-labeled root cause
  triage_date TIMESTAMP,
  created_at  TIMESTAMP DEFAULT NOW()
);

CREATE INDEX idx_project_status ON regression_runs(project, status);
CREATE INDEX idx_failure_sig ON regression_runs(failure_sig);
```

Start collecting structured data now, even before you have ML models to consume it. The most common regret among teams adopting AI for EDA is: "We wish we had been logging this data three projects ago."

CHAPTER 8

Getting Started and Future Outlook

Practical Adoption Roadmap

Adopting AI in a semiconductor design flow is not a single decision—it is a multi-phase journey that must be calibrated to your team's current capabilities, design complexity, and organizational appetite for change. The following roadmap has been refined through conversations with engineering leads at companies ranging from startups to top-10 semiconductor firms.

Phase 1: Foundation (Months 1-3)

1. **Identify one high-pain workflow.** Choose a specific, bounded problem where your team spends disproportionate effort for modest results. Good candidates: lint triage, regression failure triage, coverage closure for a specific block, or RTL code review for a new hire's code.
2. **Assess your data.** What data do you already have for this workflow? Coverage databases? Triage labels? Tool run logs? If data is sparse, begin collecting it immediately—structured logging is the highest-ROI investment you can make.
3. **Run controlled experiments.** Give three engineers access to an LLM (Claude, GPT-4, or Copilot) for RTL drafting and code review for one sprint. Measure productivity changes quantitatively: lines of code produced, time to first passing simulation, number of review iterations.

Phase 2: Targeted Deployment (Months 3-9)

1. **Deploy your first ML model.** Based on Phase 1 assessment, build or integrate one ML-powered capability. If you chose lint triage, train a classifier on historical waiver data. If you chose coverage, experiment with an RL-based stimulus agent on one block.
2. **Integrate with CI/CD.** The model must operate within your existing flow—not as a standalone experiment. This means producing output in the format your

engineers consume (e.g., annotated lint reports, prioritized test lists, coverage predictions in your dashboard).

3. **Establish feedback loops.** When engineers override the model's recommendations, capture that feedback as training data. This continuous learning loop is what separates production ML from one-off experiments.

Phase 3: Scaling (Months 9–18)

1. **Expand to additional workflows.** Apply lessons from the first deployment to other pain points. By this phase, your data infrastructure and ML tooling are mature enough to support faster iteration.
2. **Evaluate commercial AI-enhanced tools.** With internal experience, you can evaluate vendor claims critically. Run head-to-head comparisons: DSO.ai vs. Cerebrus on your designs, or commercial coverage AI vs. your internal models.
3. **Build organizational capability.** Hire or train ML engineers who understand EDA. This is a scarce skill set, but your Phase 1–2 experience gives you credibility and a concrete problem statement for recruiting.

Start Small: Pick One Workflow

The single most common failure mode in AI adoption is trying to do too much at once. A team that simultaneously tries LLM-based RTL generation, ML-based verification, and AI-driven P&R will likely deliver none of them to production quality. Start with one workflow, prove value, build credibility, and expand.

The best starting workflows share three characteristics:

- **Measurable outcome:** You can quantify improvement (coverage percentage, triage time, defect detection rate).
- **Available data:** You have historical data to train on or evaluate against.
- **Tolerable risk:** Failure of the AI component does not block the project—it just means falling back to the manual process.

Building Internal Skills

Your team needs to develop competency in three areas:

Skill Area	Who Needs It	Learning Path
ML fundamentals	All engineers (awareness level)	Andrew Ng's ML Specialization, fast.ai course, internal lunch-and-learns
Prompt engineering for HDL	RTL designers, verification engineers	Practice with LLMs on real design tasks; develop internal prompt libraries
ML engineering for EDA	1–2 dedicated engineers per team	PyTorch/TensorFlow, RL (Stable Baselines3), GNN (PyG), experiment tracking (MLflow/W&B)
Data engineering	CAD/infrastructure team	ETL pipelines, data warehousing, structured logging for EDA tool outputs

You do not need to hire a team of ML PhDs. The most effective pattern is embedding one or two ML-capable engineers within the design or verification team, working alongside domain experts. The domain experts know which problems matter; the ML engineers know which techniques are tractable.

Data Strategy

Your proprietary design data is your most valuable asset for AI adoption—far more valuable than any off-the-shelf model. A deliberate data strategy includes:

- **Structured logging:** Instrument your EDA scripts to emit structured logs (JSON, Parquet) with tool configurations, runtime metrics, and quality-of-result data. Capture this for every synthesis, P&R, and simulation run.
- **Label collection:** When engineers triage failures, waive lint warnings, or review code, capture their decisions as labels. These labeled datasets are the fuel for supervised ML.
- **Version linking:** Connect RTL versions (git SHAs) to synthesis results, simulation coverage, and tapeout outcomes. This enables training models that predict downstream outcomes from RTL features.

- **Data governance:** Classify data by sensitivity. Public benchmarks and open-source IP can be used freely; internal RTL and verification data requires access controls and may restrict which ML tools you can use (on-premises vs. cloud).

Measuring ROI

Engineering managers need to justify AI investments with concrete metrics. The most defensible ROI calculations use time savings as the primary unit:

- **Verification:** Reduction in time-to-coverage-closure (e.g., from 6 weeks to 2 weeks). Reduction in regression triage time (e.g., from 3 person-days to 0.5 person-days per milestone).
- **Physical design:** Reduction in implementation iterations (e.g., from 15 manual iterations to 3 AI-guided iterations). PPA improvement quantified in MHz, mm², and mW.
- **Code generation:** Reduction in time from specification to first compilable RTL draft. Measure carefully—include the review and correction time, not just the generation time.
- **Bug detection:** Number of bugs caught by ML that were missed by traditional tools. Ideally, measure this retroactively: would the model have caught bugs that escaped to silicon on previous projects?

KEY CONCEPT: MEASURE THE FULL LOOP

Do not measure only the AI model's output quality in isolation. Measure the total human + AI workflow time, including the time engineers spend reviewing, correcting, and integrating AI outputs. A model that generates RTL in 30 seconds but requires 4 hours of debugging is not saving time compared to 2 hours of manual coding.

Future Directions

Autonomous Chip Design

The long-term trajectory points toward increasingly autonomous design systems. Today, AI assists with individual steps: RTL drafting, parameter optimization, coverage closure. The next generation will chain these steps together, with AI agents

managing multi-step design flows end-to-end—from specification to layout, with human engineers providing oversight and strategic direction rather than hands-on execution.

Projects like NVIDIA's ChipNeMo and Google's AlphaChip are early steps in this direction. The missing pieces are reliable AI-driven verification (the hardest problem) and formal methods integration (to provide guarantees that statistical testing cannot).

AI Co-Pilots for Design Engineers

The most impactful near-term development is the evolution of AI from a batch tool to an interactive co-pilot. Imagine an AI assistant that:

- Monitors your RTL edits in real time and flags potential issues before you save the file.
- Suggests assertions based on the code you just wrote and the specification document in your project.
- Predicts the synthesis impact of a code change (area, timing) without running the tool.
- Generates targeted test stimuli for the specific logic you just modified.

This is not speculative—it is the natural extension of IDE-integrated LLMs (like GitHub Copilot) to the hardware design domain, augmented with domain-specific models for synthesis prediction, coverage estimation, and formal property generation.

Multi-Modal Models

Future AI systems will understand chip design across modalities: reading specification documents (natural language), interpreting block diagrams and schematics (vision), analyzing RTL (code), and reasoning about timing/power reports (structured data). Multi-modal models that fuse these information sources will enable higher-level design automation—for example, generating an RTL implementation directly from a microarchitecture specification document that includes both text descriptions and block diagrams.

Recommended Resources

Papers

- Mirhoseini et al., "A Graph Placement Methodology for Fast Chip Design" (Nature, 2021)—the AlphaChip paper.
- Liu et al., "ChipNeMo: Domain-Adapted LLMs for Chip Design" (NVIDIA, 2023)—the definitive work on fine-tuning LLMs for EDA.
- Thakur et al., "Benchmarking Large Language Models for Automated Verilog RTL Code Generation" (DATE, 2023)—evaluation framework for RTL-generating LLMs.
- Huang et al., "Machine Learning for Electronic Design Automation: A Survey" (ACM TODAES, 2021)—comprehensive survey of ML applications across the EDA flow.
- Wu et al., "GAMORA: Graph Learning based Symbolic Reasoning for Large-Scale Boolean Networks" (DAC, 2023)—GNN-based formal reasoning for hardware.

Courses and Tutorials

- **Stanford EE271 + CS229 combo:** Pair a core VLSI course with a core ML course to build practical ML-for-EDA foundations.
- **fast.ai:** Practical Deep Learning for Coders—the best starting point for engineers new to ML.
- **OpenROAD tutorials:** Hands-on experience with an open-source RTL-to-GDSII flow with ML components.
- **Stable Baselines3 documentation:** Practical RL implementation for coverage-driven verification experiments.

Communities

- **WOSET (Workshop on Open-Source EDA Technology):** Annual workshop at ICCAD, focusing on open-source EDA including ML-enhanced tools.
- **MLCAD Symposium:** ACM/IEEE International Symposium on Machine Learning for CAD.
- **CIRCT community:** Active development of ML-friendly hardware compiler infrastructure (GitHub: llvm/circt).
- **r/chipdesign, HardwareAI Discord:** Practitioner communities discussing AI/ML adoption in real design teams.

CHAPTER 9

AI Development Environment for VLSI Engineers

The preceding chapters focused on *what* AI can do across the chip design flow. This chapter covers *how* to set up and configure AI tools for daily VLSI work. We introduce three concepts that separate productive AI-assisted engineering from ad-hoc chatbot usage: **rules** (persistent coding standards for AI), **MCP** (a protocol that lets AI invoke EDA tools directly), and **skills** (reusable, composable AI capabilities for common workflows). We close with advanced prompt engineering techniques that go well beyond the template library in Appendix A.

9.1 Project Rules for AI-Assisted HDL Development

When you use an AI coding assistant—whether Cursor, VS Code with Copilot, or a chat-based model—it starts with no knowledge of your project's coding standards. Every session, you would need to repeat: "use synchronous reset," "non-blocking assignments in sequential blocks," "name inputs with the `i_` prefix." This is tedious, error-prone, and scales poorly across a team.

Rules solve this problem. A rule is a persistent, project-level instruction that is automatically injected into the AI's context whenever it generates or reviews code. Rules ensure that every AI interaction respects your team's conventions without manual prompting.

Rule Formats and Activation

Different AI tools support different rule mechanisms, but the concept is universal:

Tool	Rule Location	Format
Cursor IDE	<code>.cursor/rules/*.mdc</code>	Markdown with YAML frontmatter; supports glob-based auto-attach
Cursor (legacy)	<code>.cursorrules</code>	Plain text file at project root
Generic (any LLM)	<code>AGENTS.md</code>	Markdown file; works as a system prompt when pasted or referenced
ChatGPT / Claude	Custom Instructions / System Prompt	Plain text configured in the UI or API

Modern IDEs like Cursor support four activation modes:

- **Always Apply:** Active on every interaction. Use for universal coding standards.
- **Auto-Attached:** Triggered when the file being edited matches a glob pattern (e.g., `**/*.sv`). Use for language-specific rules.
- **Agent-Decided:** The AI reads the rule's description and decides whether it is relevant. Use for workflow-specific guidance (e.g., "apply when generating testbenches").
- **Manual:** Applied only when explicitly referenced with `@rule-name`. Use for specialized, infrequent tasks.

Rule Template 1: RTL Coding Standard

This rule enforces synthesizable SystemVerilog conventions. Save it as `.cursor/rules/rtl-coding-standard.mdc` with globs: `["**/*.sv", "**/*.v"]`:


```

---
description: RTL coding standard for all SystemVerilog and Verilog files
globs: ["**/*.sv", "**/*.v"]
alwaysApply: false
---
```

RTL Coding Standard

You are generating synthesizable RTL. Follow these rules strictly:

Reset and Clocking

- Use synchronous, active-low reset (rst_n) unless the spec explicitly requires otherwise.
- All sequential logic uses ``always_ff @(posedge clk)``.
- All combinational logic uses ``always_comb``.
- Never mix blocking and non-blocking in the same always block.

Assignments

- Non-blocking (``<=``) in ``always_ff`` blocks.
- Blocking (``=``) in ``always_comb`` blocks.
- No assignments to the same signal from multiple always blocks.

Naming Conventions

- Inputs: ``i_`` prefix (e.g., ``i_data``, ``i_valid``)
- Outputs: ``o_`` prefix (e.g., ``o_ready``, ``o_result``)
- Internal wires: ``w_`` prefix
- Internal registers: ``r_`` prefix
- Parameters: ``UPPER_SNAKE_CASE``
- Module names: ``lower_snake_case``

Synthesis Subset

- No ``initial`` blocks (use reset instead).
- No ``#delay`` constructs.
- Every ``if`` in combinational logic must have an ``else``.
- Every ``case`` must have a ``default``.
- No inferred latches—ever.

Structure

- Separate FSM control logic from datapath.
- One module per file. Filename matches module name.
- Define all register reset values explicitly.

Rule Template 2: UVM Verification

Save as `.cursor/rules/uvm-verification.mdc` with `globs: ["**/tb/**/*.sv", "**/verif/**/*.sv"]`:

```

---
description: UVM verification conventions for testbench files
globs: ["**/tb/**/*.sv", "**/verif/**/*.sv"]
alwaysApply: false
---

# UVM Verification Standard

You are generating UVM testbench code. Follow these conventions:

## UVM Macros and Registration
- Use `uvm_component_utils` / `uvm_object_utils` for all classes.
- Always call `super.new(name, parent)` in constructors.
- Use factory overrides, not direct construction.

## Sequence and Transaction Style
- Transactions extend `uvm_sequence_item`.
- Sequences extend `uvm_sequence #(transaction_type)`.
- Use `start_item()` / `finish_item()` pattern.
- Randomize with inline constraints where possible.

## Assertions and Coverage
- Name every assertion with a label (e.g., `assert_no_overflow`).
- Include `$error()` messages with context in every assertion.
- Define `covergroup` inside the monitor or a dedicated collector.
- Cross-cover related signals when state space is manageable.

## Scoreboard
- Use TLM analysis ports for monitor-to-scoreboard communication.
- Log mismatches with `uvm_error` including expected vs. actual values.

```

Rule Template 3: Tcl/SDC Constraints

Save as `.cursor/rules/tcl-sdc.mdc` with `globs: ["**/*.tcl", "**/*.sdc"]`:

```

---
description: Tcl and SDC constraint file conventions
globs: ["**/*.tcl", "**/*.sdc"]
alwaysApply: false
---

# Tcl/SDC Conventions

## SDC Constraints
- Define all clocks with `create_clock` before any derived constraints.
- Use `create_generated_clock` for PLL/divider outputs.
- Group false paths and multicycle paths with comments explaining rationale.
- Set input/output delays relative to the correct clock domain.

## Tcl Scripting
- Use `proc` for reusable operations.
- Prefer `get_pins`/`get_ports`/`get_cells` with explicit `-hierarchical` w
- Quote all collection arguments to avoid glob expansion surprises.
- Add `puts` progress messages for long-running scripts.

## File Organization
- Separate timing constraints (.sdc) from tool-specific scripts (.tcl).
- Version-control all constraint files alongside RTL.

```

KEY CONCEPT: RULES ARE LIVING DOCUMENTS

Rules should evolve with your project. Start with a minimal set (reset style, naming, synthesis subset), then add rules as you discover recurring AI mistakes. Review your rules at each project milestone, the same way you review your coding guidelines.

9.2 MCP: Model Context Protocol for EDA

Rules tell the AI *how* to write code. MCP lets the AI *run* tools. The **Model Context Protocol** (MCP), released as an open standard by Anthropic in late 2024, is a JSON-RPC-based protocol that allows AI assistants to invoke external tools, read files, and query databases. Think of it as a standardized API layer between an LLM and the outside world.

For VLSI engineers, MCP is transformative because it closes the loop between code generation and code validation. Instead of the engineer manually copying generated RTL into a terminal, running a lint check, copying the errors back into the chat, and asking for a fix, an MCP-enabled AI can do this autonomously.

How MCP Works

The architecture is client-server:

```
graph LR
  A[AI Assistant] -->|"JSON-RPC"| B[MCP Server]
  B -->|"exec"| C[Yosys]
  B -->|"exec"| D[Icarus Verilog]
  B -->|"exec"| E[Verilator]
  B -->|"read"| F[Design Files]
  C -->|"results"| B
  D -->|"results"| B
  E -->|"results"| B
  B -->|"response"| A
```

Figure 9.1: MCP Architecture for EDA Tool Integration

1. **MCP Client** (the AI assistant, e.g., Cursor or Claude Desktop) discovers available tools by querying the MCP server.
2. **MCP Server** exposes EDA tools as callable functions with typed parameters and return values. Each tool has a name, description, and JSON schema for its inputs.
3. When the AI decides it needs to compile a module, it sends a `tools/call` request to the server with the appropriate parameters.
4. The server executes the EDA tool, captures stdout/stderr, and returns structured results to the AI.
5. The AI interprets the results (compilation errors, lint warnings, simulation output) and takes the next action—fixing the code, adjusting constraints, or reporting to the engineer.

MCP4EDA: A Concrete Implementation

Production Practice (emerging): MCP4EDA and AutoEDA are open-source research implementations. Production deployment requires validation on your tool versions and design complexity.

MCP4EDA [27] is the first published MCP server for EDA, integrating Yosys (synthesis), Icarus Verilog (simulation), OpenLane (place-and-route), and GTKWave (waveform analysis). It demonstrates a closed-loop flow where an AI agent:

1. Generates RTL from a natural-language specification.

2. Synthesizes with Yosys and reads back area/timing reports.
3. Simulates with Icarus Verilog and checks for assertion failures.
4. Iterates on the RTL based on synthesis and simulation feedback.
5. Runs place-and-route through OpenLane and reports final PPA metrics.

Published results show 15-30% timing improvements and 10-20% area reduction compared to default synthesis flows, achieved through iterative parameter exploration guided by actual back-end results rather than front-end estimates.

Example: Minimal MCP Server for Yosys + Icarus Verilog

The following shows the structure of an MCP tool definition for synthesis. In practice you would use a framework like the MCP Python SDK or TypeScript SDK to expose these tools:

```

# MCP tool definitions for a basic open-source EDA server
# Each tool becomes callable by the AI assistant

tools = [
    {
        "name": "yosys_synth",
        "description": "Synthesize a SystemVerilog file with Yosys. Returns",
        "inputSchema": {
            "type": "object",
            "properties": {
                "file_path": {"type": "string", "description": "Path to the"},
                "top_module": {"type": "string", "description": "Top-level"},
                "target": {"type": "string", "description": "Target library"}
            },
            "required": ["file_path", "top_module"]
        }
    },
    {
        "name": "iverilog_sim",
        "description": "Compile and simulate a Verilog testbench with Icarus",
        "inputSchema": {
            "type": "object",
            "properties": {
                "files": {"type": "array", "items": {"type": "string"}},
                "top_module": {"type": "string"}
            },
            "required": ["files"]
        }
    },
    {
        "name": "verilator_lint",
        "description": "Run Verilator in lint-only mode. Returns warnings a",
        "inputSchema": {
            "type": "object",
            "properties": {
                "file_path": {"type": "string"},
                "top_module": {"type": "string"}
            },
            "required": ["file_path", "top_module"]
        }
    }
]

```

The Closed-Loop Workflow

With MCP, the RTL development loop from Chapter 3 becomes fully automated between generation and verification:

```
graph TD
    A[Engineer: Write Spec] --> B[AI: Generate RTL]
    B -->|"MCP: verilog_lint"| C{Lint Clean?}
    C -->|No| B
    C -->|Yes| D[AI: Generate Testbench]
    D -->|"MCP: iverilog_sim"| E{Tests Pass?}
    E -->|No| B
    E -->|Yes| F[AI: Synthesize]
    F -->|"MCP: yosys_synth"| G[Report PPA to Engineer]
    G --> H[Engineer: Review and Commit]
```

Figure 9.2: MCP-Enabled Closed-Loop RTL Development

The engineer's role shifts from executing tools manually to reviewing the AI's work and making architectural decisions. The AI handles the compile-fix-recompile iteration cycle that consumes hours of engineering time on routine blocks.

Security and IP Considerations

MCP raises important security questions for semiconductor companies:

- **On-premises only for proprietary designs.** MCP servers that invoke commercial EDA tools (Design Compiler, VCS, JasperGold) must run on your internal infrastructure. Never route proprietary RTL through cloud-hosted MCP servers.
- **Tool license compliance.** Automated MCP-driven tool invocations consume the same licenses as manual runs. Ensure your license pool can handle increased utilization from AI-driven iteration.
- **Audit logging.** Log every MCP tool invocation (tool name, parameters, timestamp, user) for traceability. This is essential for IP governance and for debugging AI-generated designs post-tapeout.
- **Sandboxing.** Run MCP servers in containers with restricted filesystem access. The AI should only be able to read and write within the project workspace, not traverse the broader design repository.

9.3 AI Agent Skills for EDA Workflows

A **skill** is a reusable, self-contained AI capability that packages domain knowledge, tool invocations, and output formatting into a single callable unit. If rules are "how to write code" and MCP is "how to run tools," skills are "how to accomplish a complete task."

Skills compose rules and MCP into end-to-end workflows. A well-designed skill takes a structured input (e.g., a module file path), executes a sequence of actions (lint, simulate, analyze), and produces a structured output (a review report, a coverage summary, a set of suggested fixes).

Skill Architecture

Component	Purpose	Example
Input specification	What the skill needs to start	Path to a .sv file and the module's spec document
Rules context	Coding standards the skill respects	RTL coding standard rule (Section 9.1)
Tool chain (MCP)	EDA tools the skill invokes	Verilator lint, Icarus Verilog sim, Yosys synth
Reasoning template	Step-by-step logic the AI follows	"First lint, then check resets, then verify FSM completeness"
Output format	Structured result for the engineer	Markdown report with severity-ranked findings

Skill 1: RTL Review

This skill performs a comprehensive code review that combines static analysis with AI reasoning:

Skill: RTL Code Review

Trigger: Engineer requests review of a SystemVerilog module

Inputs

- `file_path`: Path to the .sv file under review
- `spec_doc` (optional): Path to the design specification

Procedure

1. Read the source file and the project's RTL coding standard rule.
2. Run `verilator_lint` via MCP. Collect all warnings and errors.
3. Analyze the code for:
 - Reset completeness: every register has an explicit reset value.
 - FSM safety: all states have defined transitions; no deadlocks.
 - CDC risk: signals crossing clock domain boundaries without synchronize.
 - Latch inference: incomplete if/case in combinational blocks.
 - Naming violations: signals not following i_/o_/w_/r_ conventions.
4. If `spec_doc` is provided, cross-reference the implementation against the specification for missing functionality or interface mismatches.
5. Produce a structured report:
 - CRITICAL: issues that will cause functional bugs
 - WARNING: issues that may cause synthesis or timing problems
 - STYLE: coding standard violations
 - SUGGESTION: improvements for readability or maintainability

Output Format

Markdown report with one section per finding, each containing:
line number, severity, description, and suggested fix.

Skill 2: Coverage Gap Analysis

Skill: Coverage Gap Analysis

Trigger: Engineer provides a coverage database after regression

Inputs

- `coverage_report`: Path to the coverage report (UCDB, HTML, or text)
- `design_spec` (optional): Spec document for context

Procedure

1. Parse the coverage report to identify uncovered bins and low-hit cross p
2. Group uncovered points by coverage group and rank by verification priori
3. For each uncovered group, analyze why it might be hard to reach:
 - Is the scenario unreachable (dead code)?
 - Does it require a specific multi-cycle input sequence?
 - Is it blocked by a constraint that is too restrictive?
4. Suggest targeted test strategies:
 - Specific constraint modifications
 - Directed test sequences
 - Assertion-based coverage points to add
5. Estimate the effort to close each gap (low/medium/high).

Output Format

Table: coverage group | uncovered bins | root cause hypothesis |
suggested action | effort estimate

Skill 3: Timing Closure Assistant

Skill: Timing Closure Assistant

Trigger: Engineer provides an STA timing report with violations

Inputs

- `timing_report`: Path to STA report (e.g., PrimeTime or OpenSTA output)
- `rtl_sources` (optional): Paths to relevant RTL files

Procedure

1. Parse the timing report to extract all violating paths.
2. Classify violations: setup vs. hold, clock domain, path type.
3. For each critical path:
 - Identify the logic depth and the bottleneck stages.
 - Check if the path crosses a clock domain boundary.
 - Analyze whether the violation is in control logic or datapath.
4. Suggest fixes ranked by impact and invasiveness:
 - Constraint fix: is the path a false path or multicycle path?
 - RTL restructuring: pipeline insertion, logic balancing, retiming.
 - Floorplan adjustment: macro placement or partition boundary change.
5. If RTL sources are provided, generate specific code changes.

Output Format

Ranked list of violations with: slack, path summary, root cause, and recommended fix (constraint change or RTL diff).

Skill 4: Testbench Scaffold Generator

Skill: Testbench Scaffold Generator

Trigger: Engineer provides a module interface for testbench creation

Inputs

- `module_file`: Path to the DUT's .sv file
- `framework`: "cocotb" | "uvm" | "basic_sv"
- `protocol` (optional): "axi4", "axi_stream", "wishbone", "custom"

Procedure

1. Parse the module to extract the port list, parameters, and clock/reset signals
2. Apply the project's verification rule (Section 9.1) for style convention
3. Generate the testbench scaffold:
 - Clock and reset generation
 - DUT instantiation with all ports connected
 - Driver: randomized stimulus with protocol-aware constraints
 - Monitor: output capture and protocol checking
 - Scoreboard: basic expected-vs-actual comparison
 - Initial test: reset sequence + 100 random transactions + drain
4. If a protocol is specified, use protocol-specific bus functional models (e.g., cocotbext-axi for AXI, or UVM agent for standard protocols).
5. Run `iverilog_sim` via MCP to verify the scaffold compiles and the reset sequence executes without errors.

Output Format

Complete testbench file(s) ready to run, plus a summary of what was generated and what the engineer should customize.

KEY CONCEPT: SKILLS COMPOSE

The real power emerges when skills chain together. An engineer says "review this module and generate a testbench for it." The AI invokes the RTL Review skill first, fixes any critical issues it finds, then invokes the Testbench Scaffold skill on the corrected code. Each skill respects the same project rules and uses the same MCP tools, ensuring consistency across the entire workflow.

9.4 Advanced Prompt Engineering for EDA Workflows

The prompt templates in Appendix A are effective starting points, but production-quality AI-assisted design requires more sophisticated prompting techniques. This section covers the methods that consistently produce the best results in VLSI workflows.

Chain-of-Thought for RTL Design

Chain-of-thought (CoT) prompting asks the model to reason through the design step-by-step before writing code. This dramatically reduces errors on complex logic because the model "thinks through" edge cases before committing to an implementation.

Design a round-robin arbiter with 4 requestors. Before writing any code, work through the design step by step:

1. **Interface definition:** List all ports with widths and directions.
2. **State analysis:** What states does the FSM need?
What are the transition conditions?
3. **Priority logic:** How does the round-robin pointer advance?
What happens when no requestors are active?
4. **Edge cases:** What if multiple requestors assert simultaneously?
What if a granted requestor de-asserts before being serviced?
5. **Reset behavior:** What is the initial state? Initial priority pointer?

After completing the analysis, write the SystemVerilog implementation.

CoT is especially valuable for FSMs, protocol implementations, and any module with non-trivial state transitions. The step-by-step reasoning also makes it easier for the reviewing engineer to verify the AI's design intent before examining the code.

Few-Shot Prompting with Golden Examples

Including one or two reviewed, correct code snippets in the prompt gives the model a concrete reference for style, structure, and quality. This is more effective than describing conventions in text because the model pattern-matches against the example.

Here is an example of a well-written FIFO module that follows our coding st

```

```systemverilog
module sync_fifo #(
 parameter int DATA_W = 8,
 parameter int DEPTH = 16
) (
 input logic i_clk,
 input logic i_rst_n,
 input logic [DATA_W-1:0] i_wdata,
 input logic i_wr_en,
 input logic i_rd_en,
 output logic [DATA_W-1:0] o_rdata,
 output logic o_full,
 output logic o_empty
);
 // ... implementation follows our standard ...
endmodule
```

```

Now write a dual-clock asynchronous FIFO with the same coding style, conventions, and quality level. Use gray-code pointers for CDC.

Tool-in-the-Loop Prompting

When using MCP or manual copy-paste workflows, structuring the error feedback as a prompt continuation produces better fixes than starting a new conversation:

The module you generated has the following Verilator lint errors:

```

```
%Warning-LATCH: fifo_ctrl.sv:47: Latch inferred for signal
 'fifo_ctrl.wr_ptr_next' (not all control paths of combinational
 block assign a value)
```

```
%Warning-WIDTH: fifo_ctrl.sv:62: Operator ASSIGN expects 5 bits
 on the Assign RHS, but Assign RHS's SUB generates 32 bits.
```

```

Fix these specific issues:

1. The LATCH warning: ensure all paths in the combinational block assign `wr_ptr_next`.
2. The WIDTH warning: explicitly size the subtraction result to match the LHS width.

Return **ONLY** the corrected module. Do not change anything else.

Retrieval-Augmented Prompting

For designs that must conform to a specification document, attaching relevant spec sections to the prompt grounds the AI's output in your actual requirements rather than generic training data:

Context: The following is Section 3.2 of our DMA Controller specification.

[PASTE SPEC SECTION: register map, transfer modes, interrupt behavior]

Based on this specification, generate the register bank module for the DMA controller. Implement all registers described in the spec with the exact field definitions, reset values, and access types (RW, RO, W1C) specified above.

Flag any ambiguities in the spec as comments in the generated code.

This technique is particularly powerful when combined with rules (which set the coding style) and MCP (which validates the generated code against the spec by running simulation).

Structured Output Prompting

For analysis tasks (timing review, coverage analysis, code review), requesting structured output makes the AI's response directly actionable:

Analyze this timing report and respond in EXACTLY this format:

Violation Summary

| # | Slack (ns) | From | To | Domain | Type |
|---|------------|------|-----|--------|------------|
| 1 | ... | ... | ... | ... | setup/hold |

Root Cause Analysis

For each violation:

- **Path #N**: [one-sentence root cause]
 - Logic depth: [N levels]
 - Bottleneck: [specific gate/operation]
 - Fix category: [constraint / RTL / floorplan]

Recommended Actions (ranked by impact)

1. [Most impactful fix first]
2. ...

Do not include any other commentary.

The combination of rules, MCP, skills, and advanced prompting transforms AI from a code-suggestion tool into an integrated engineering assistant. Start with rules (low effort, immediate payoff), add MCP when you want automated validation, and define skills as your team's AI workflows mature.

Appendix A: Prompt Library

Purpose: Copy-paste templates for common workflows. Replace [BRACKETED_TEXT] with your project specifics.

A.1 RTL Generation Template

You are an expert RTL designer. Write a synthesizable SystemVerilog module

1. ****Interface:****

- Clock: [CLK_NAME] (posedge)
- Reset: [RST_NAME] ([synchronous/asynchronous], [active-high/active-low])
- Inputs: [LIST_INPUTS_AND_WIDTHS]
- Outputs: [LIST_OUTPUTS_AND_WIDTHS]

2. ****Functionality:****

- [DESCRIBE_CORE_LOGIC_STEP_BY_STEP]
- [DESCRIBE_CORNER_CASES]

3. ****Constraints:****

- Use [blocking/non-blocking] assignments correctly.
- Separate control logic (FSM) from datapath.
- Ensure reset values are defined for all registers.
- Use clear, descriptive signal names.

4. ****Output Format:****

- Provide ONLY the SystemVerilog code.
- Add comments explaining complex logic.

A.2 Code Review / Bug Hunting Template

Act as a senior Verification Engineer. Review the following SystemVerilog code:

Focus specifically on:

1. **Clock Domain Crossings (CDC):** Are signals crossing domains without synchronization?
2. **Latches:** Are there incomplete case/if statements?
3. **Reset Logic:** Are resets handled consistently?
4. **FSM Safety:** Are there potential deadlocks or unreachable states?
5. **Width Mismatches:** Are there implicit truncations or extensions?

Code to review:

```
```systemverilog
[PASTE_CODE_HERE]
```
```

Output a bulleted list of issues. For each issue, cite the line number, expected behavior, and the actual behavior.

A.3 SystemVerilog Assertion (SVA) Generation

Write SystemVerilog Assertions (SVA) for the following protocol requirement:

"Requirement: [PASTE_REQUIREMENT_TEXT_FROM_SPEC]"

Context:

- Clock: [CLK_NAME]
- Reset: [RST_NAME]
- Signal names: [LIST_RELEVANT_SIGNALS]

Instructions:

1. Define a property that captures this requirement.
2. assert the property with a meaningful error message.
3. cover the property to ensure the scenario is actually stimulated.
4. Use standard SVA operators (\rightarrow , \Rightarrow , ##N, \$rose, \$fell).

A.4 Testbench Scaffolding (UVM/Cocotb)

Create a [UVM/Cocotb] testbench skeleton for a module with this interface:

[PASTE_MODULE_INTERFACE]

The testbench should include:

1. A driver to drive random inputs.
2. A monitor to capture outputs.
3. A scoreboard to check protocol correctness (basic sanity checks).
4. A test that runs for [NUMBER] cycles/transactions.

Use modern [SystemVerilog/Python] best practices.

A.5 Tcl/SDC Constraint Generation

Generate SDC timing constraints for the following block:

Block name: [BLOCK_NAME]

Primary clock: [CLK_NAME], [FREQUENCY] MHz, [DUTY_CYCLE]% duty cycle

Generated clocks: [LIST_OF_DIVIDED_OR_PLL_CLOCKS_WITH_RATIOS]

Interface timing:

- Input ports: [LIST_PORTS] arrive [N] ns after [CLK_NAME] rising edge
- Output ports: [LIST_PORTS] must be stable [N] ns before [CLK_NAME] rising edge

Clock domain crossings:

- [SRC_CLK] to [DST_CLK]: [synchronizer type, e.g., 2FF / handshake / async]

Known false paths:

- [DESCRIPTION_OF_FALSE_PATHS, e.g., "test_mode scan signals to functional"]

Instructions:

1. Define all clocks with ``create_clock`` and ``create_generated_clock``.
2. Set ``set_input_delay`` and ``set_output_delay`` for all I/O ports.
3. Add ``set_false_path`` for documented false paths and CDC crossings that use proper synchronization.
4. Add ``set_multicycle_path`` where specified.
5. Add comments explaining every constraint group.

A.6 Timing Report Analysis

You are a senior physical design engineer. Analyze the following STA timing and provide actionable recommendations.

```

[PASTE\_TIMING\_REPORT\_EXCERPT: violating paths with slack, logic levels, cell  
```

For each violating path, provide:

1. **Path summary**: source register -> combinational stages -> destination
2. **Bottleneck identification**: which cell or net contributes the most delay
3. **Root cause category**: logic depth / high-fanout net / long wire / clock incorrect constraint.
4. **Recommended fix** (ranked by least invasive first):
 - a. Constraint fix (false path, multicycle, or clock adjustment)
 - b. Synthesis directive (size_only, dont_touch removal, effort level)
 - c. RTL change (pipelining, logic restructuring, retiming)
 - d. Floorplan change (macro placement, partition boundary)
5. **Estimated slack recovery** for each recommended fix.

Output as a structured table.

A.7 Verification Plan Generation

Generate a verification plan for the following design block:

Block name: [BLOCK_NAME]

Specification: [PASTE_OR_SUMMARIZE_KEY_SPEC_SECTIONS]

Interfaces: [LIST: protocol, width, direction]

Key features to verify:

- [FEATURE_1]
- [FEATURE_2]
- [FEATURE_N]

The verification plan should include:

1. **Feature extraction**: List every testable feature from the spec.
2. **Coverage model**:
 - Functional coverage groups with bin definitions.
 - Cross-coverage points for interacting features.
 - Identify bins that are likely unreachable and should be excluded.
3. **Test strategy** for each feature:
 - Directed tests (specific scenarios)
 - Constrained-random approach (constraints and distributions)
 - Assertion-based checks (SVA properties)
4. **Corner cases**: List edge conditions that require targeted stimulus.
5. **Coverage closure criteria**: Minimum coverage thresholds for sign-off.
6. **Estimated effort**: T-shirt sizing (S/M/L) for each test category.

Format as a structured markdown document suitable for a design review.

A.8 FPGA-to-ASIC Migration Review

Review the following RTL module that was originally written for FPGA and flag all issues that need to be addressed for ASIC implementation.

```
```systemverilog
[PASTE_MODULE_CODE]
```
```

Target ASIC process: [PROCESS_NODE, e.g., TSMC 7nm]

Target library: [STANDARD_CELL_LIBRARY]

Check for and report:

1. **FPGA-specific constructs**: Block RAM inference attributes, DSP pragma, FPGA primitive instantiations (BUFG, IBUF, PLL, MMCM, etc.).
2. **Clock and reset**: FPGA-style clock management (BUFG, MMCM) that must be replaced with ASIC PLL/clock-tree elements.
3. **Memory inference**: Inferred block RAMs that should become compiled SRAM with proper BIST wrappers.
4. **Synthesis subset**: Any constructs that are supported by FPGA synthesis but not by ASIC synthesis (e.g., `initial` blocks for register initialization).
5. **Timing assumptions**: Hard-coded timing values that assume FPGA clock.
6. **DFT readiness**: Missing scan chain support, test mode signals, or BIST.
7. **Power considerations**: Missing clock gating opportunities, always-on logic that should be power-gated.

For each issue, provide: location, severity (blocking/warning/info), description, and the specific ASIC-compatible replacement.

Appendix B: Review Checklists

Purpose: Standardize human-in-the-loop review for AI-generated outputs.

- **RTL checklist:** Interface/spec alignment, synthesis legality, CDC/reset sanity, and simulation behavior.
- **SVA checklist:** Correct temporal operators, no vacuity, proper scoping, and meaningful failure messages.
- **Verification checklist:** Coverage impact, false-positive rate, and regression stability across seeds.

Appendix C: Metrics Dashboard Template

Purpose: Track pilot impact with weekly metrics.

- **Cycle-time metrics:** Time to first compile, time to first passing sim, time to coverage closure.
- **Quality metrics:** Bugs found pre-merge, escaped defects, assertion pass/fail trends.
- **Efficiency metrics:** Regression triage hours, lint warning reduction, compute cost per milestone.
- **Adoption metrics:** % changes using AI assistance, engineer override rate, acceptance/rework ratio.

Appendix D: Data Governance and IP Policy

Purpose: Define safe AI usage boundaries for proprietary design data.

- **Data classification:** Public/open IP, internal non-sensitive, restricted proprietary, export-controlled.
- **Deployment policy:** Which data classes are allowed in cloud APIs versus on-prem inference.
- **Auditability:** Log prompts, model versions, and review outcomes for traceability.
- **Approval workflow:** Security/legal sign-off requirements for new tools and integrations.

Appendix E: AI Rules Reference for VLSI Projects

Purpose: Copy-paste rule files for configuring AI coding assistants in VLSI projects. These rules work with Cursor IDE (`.cursor/rules/` directory), and the content can be adapted as system prompts for any LLM.

Setup Instructions

1. Create a `.cursor/rules/` directory in your project root.
2. Save each rule below as a separate `.mdc` file in that directory.
3. Rules with `globs` activate automatically when editing matching files.

4. For non-Cursor tools (ChatGPT, Claude, Copilot Chat), paste the rule content into your system prompt or custom instructions.

E.1 SystemVerilog RTL Rule

File: `.cursor/rules/rtl-standard.mdc`


```

---
description: Synthesizable RTL coding standard
globs: ["**/*.sv", "**/*.v", "**/*.svh"]
alwaysApply: false
---

# RTL Coding Standard

## General
- Generate ONLY synthesizable SystemVerilog (IEEE 1800-2017 synthesis subse
- No `initial` blocks, no `#delay`, no `$display` in RTL (these belong in t
- One module per file. Filename must match module name.

## Clocking and Reset
- Default: synchronous active-low reset (`i_rst_n`), posedge clock (`i_clk`
- All `always_ff` blocks: `always_ff @(posedge i_clk)`
- All registers must have explicit reset values.

## Assignments
- `always_ff`: non-blocking (`<=`) only.
- `always_comb`: blocking (`=`) only.
- Never assign the same signal from multiple always blocks.

## Naming
- Inputs: `i_` prefix. Outputs: `o_` prefix.
- Wires: `w_` prefix. Registers: `r_` prefix.
- Parameters/localparams: `UPPER_SNAKE_CASE`.
- Module names: `lower_snake_case`.

## Combinational Safety
- Every `if` must have an `else` in `always_comb`.
- Every `case` must have a `default`.
- Zero inferred latches.

## Structure
- Separate control (FSM) from datapath.
- Parameterize widths and depths—no magic numbers.
- Keep modules under 300 lines; refactor if larger.

```

E.2 UVM Verification Rule

File: `.cursor/rules/uvm-verification.mdc`

```

---
description: UVM testbench coding conventions
globs: ["**/tb/**/*.sv", "**/verif/**/*.sv", "**/test/**/*.sv"]
alwaysApply: false
---

# UVM Verification Standard

## Class Registration
- All components: `uvm_component_utils(class_name)`.
- All objects/transactions: `uvm_object_utils(class_name)`.
- Always call `super.new(name, parent)` in constructors.
- Use factory overrides; never construct with `new` directly for overridabl

## Transactions
- Extend `uvm_sequence_item`.
- Use `rand` for all stimulus fields.
- Implement `do_compare`, `do_copy`, `convert2string` for debug.

## Sequences
- Use `start_item` / `finish_item` pattern.
- Prefer inline constraints: `item.randomize()` with { ... }`.
- Keep sequences short and composable; build complex scenarios by chaining.

## Assertions
- Label every assertion: `assert_name: assert property (...)`.
- Include `$error` with formatted message showing relevant signal values.
- Pair every `assert` with a `cover` to confirm scenario reachability.

## Coverage
- Define `covergroup` in the monitor or a dedicated collector component.
- Use `type_option.weight` to prioritize critical groups.
- Bin ranges should match architectural boundaries (e.g., FIFO depth thresh

## Scoreboard
- Use TLM `analysis_port` / `analysis_export` for monitor-to-scoreboard.
- Log mismatches with `uvm_error` including expected vs. actual.
- Track transaction count and report in `report_phase`.

```

E.3 Tcl/SDC Constraints Rule

File: `.cursor/rules/tcl-sdc.mdc`

```

---
description: Tcl scripting and SDC constraint conventions
globs: ["**/*.tcl", "**/*.sdc"]
alwaysApply: false
---

# Tcl/SDC Conventions

## SDC Constraints
- Define all `create_clock` before derived constraints.
- Use `create_generated_clock` for PLL/divider outputs with `-source` and `
- Group related constraints with section comments.
- Explain every `set_false_path` and `set_multicycle_path` with a comment
  stating WHY the path is false/multicycle, not just WHAT it is.
- Set `set_input_delay` and `set_output_delay` for all I/O ports.

## Tcl Scripting
- Use `proc` for reusable operations.
- Prefer `get_*` commands (`get_pins`, `get_ports`, `get_cells`) with
  `-hierarchical` when needed.
- Quote collection arguments to avoid glob expansion.
- Use `foreach_in_collection` for iterating Synopsys collections.
- Add `puts` progress messages for scripts that run longer than 30 seconds.

## File Organization
- Separate timing constraints (.sdc) from tool-specific flow scripts (.tcl)
- Name constraint files to indicate scope: `block_name.sdc`, `top_level.sdc
- Version-control all constraint files alongside RTL.

```

E.4 Cocotb Testbench Rule

File: `.cursor/rules/cocotb-testbench.mdc`

```

---
description: Cocotb Python testbench conventions
globs: ["**/test_*.py", "**/tb_*.py", "**/cocotb/**/*.py"]
alwaysApply: false
---

# Cocotb Testbench Standard

## Test Structure
- Use `@cocotb.test()` decorator for all test functions.
- First action in every test: create clock and assert reset.
- Use `await ClockCycles(dut.clk, N)` for timing, never `await Timer(ns)`.
- Each test should be self-contained (own reset, own stimulus, own checks).

## Stimulus
- Use cocotb bus extensions (cocotbext-axi, cocotbext-spi) for standard protocols.
- Use Python `random` with a fixed seed for reproducibility.
- Log the seed at test start: `cocotb.log.info(f"Seed: {seed}")`.

## Checking
- Use Python `assert` with descriptive messages.
- Compare against a Python reference model, not hardcoded expected values.
- Check at transaction boundaries, not every clock cycle.

## Organization
- One test file per DUT module.
- Common utilities (reset, clock, bus helpers) in a shared `tb_utils.py`.
- Use `Makefile` or `pytest` runner with `SIM`, `TOPLEVEL`, `MODULE` variables.

```

References and Citation Notes

Inline markers point to representative public sources. Vendor-reported figures should be treated as directional and validated on your own designs.

1. [1] Apple newsroom, "Apple reveals M3 Ultra..." (2025):
[apple.com/newsroom/2025/03/apple-reveals-m3-ultra-taking-apple-silicon-to-a-new-extreme/](https://www.apple.com/newsroom/2025/03/apple-reveals-m3-ultra-taking-apple-silicon-to-a-new-extreme/) (https://www.apple.com/newsroom/2025/03/apple-reveals-m3-ultra-taking-apple-silicon-to-a-new-extreme/)
2. [2] NVIDIA Blackwell architecture page: [nvidia.com/en-us/data-center/technologies/blackwell-architecture/](https://www.nvidia.com/en-us/data-center/technologies/blackwell-architecture/) (https://www.nvidia.com/en-us/data-center/technologies/blackwell-architecture/)

3. [3] Wilson Research Group / Siemens EDA functional verification trend report portal: verificationacademy.com/.../2024-...-functional-verification-trend-report (<https://verificationacademy.com/topics/planning-measurement-and-analysis/wrg-industry-data-and-trends/2024-siemens-eda-and-wilson-research-group-ic-asic-functional-verification-trend-report>)
4. [4] Siemens Verification Horizons, Wilson Research Group study overview (2020): blogs.sw.siemens.com/verificationhorizons/2020/10/27/... (<https://blogs.sw.siemens.com/verificationhorizons/2020/10/27/prologue-the-2020-wilson-research-group-functional-verification-study/>)
5. [5] Semiconductor Engineering, "How Much Will That Chip Cost?": semiengineering.com/how-much-will-that-chip-cost/ (<https://semiengineering.com/how-much-will-that-chip-cost/>)
6. [6] Synopsys.ai overview: [synopsys.com/ai.html](https://www.synopsys.com/ai.html) (<https://www.synopsys.com/ai.html>)
7. [7] Cadence Cerebrus AI Studio: [cadence.com/.../cadence-cerebrus-ai-studio.html](https://www.cadence.com/en_US/home/tools/digital-design-and-signoff/soc-implementation-and-floorplanning/cadence-cerebrus-ai-studio.html) (https://www.cadence.com/en_US/home/tools/digital-design-and-signoff/soc-implementation-and-floorplanning/cadence-cerebrus-ai-studio.html)
8. [8] Siemens Solido overview/news: blogs.sw.siemens.com/thought-leadership/.../reshaping-the-ic-validation-and-characterization-industry-with-ai/ (<https://blogs.sw.siemens.com/thought-leadership/2023/02/14/reshaping-the-ic-validation-and-characterization-industry-with-ai/>)
9. [9] EE Times, "Is verification really 70 percent?": [eetimes.com/is-verification-really-70-percent/](https://www.eetimes.com/is-verification-really-70-percent/) (<https://www.eetimes.com/is-verification-really-70-percent/>)
10. [10] NVIDIA Research + arXiv, "ChipNeMo: Domain-Adapted LLMs for Chip Design": arxiv.org/abs/2311.00176 (<https://arxiv.org/abs/2311.00176>)
11. [11] Cadence Verisium AI-Driven Verification Platform: [cadence.com/.../ai-driven-verification.html](https://www.cadence.com/en_US/home/tools/system-design-and-verification/ai-driven-verification.html) (https://www.cadence.com/en_US/home/tools/system-design-and-verification/ai-driven-verification.html)
12. [12] Synopsys VC Formal overview: [synopsys.com/verification/static-and-formal-verification/vc-formal.html](https://www.synopsys.com/verification/static-and-formal-verification/vc-formal.html) (<https://www.synopsys.com/verification/static-and-formal-verification/vc-formal.html>)
13. [13] Mirhoseini et al., Nature 2021, "A graph placement methodology for fast chip design": [nature.com/articles/s41586-021-03544-w](https://www.nature.com/articles/s41586-021-03544-w) (<https://www.nature.com/articles/s41586-021-03544-w>)
14. [14] Selsam et al., "Learning a SAT Solver from Single-Bit Supervision" (NeuroSAT): arxiv.org/abs/1802.03685 (<https://arxiv.org/abs/1802.03685>)
15. [15] Synopsys press release (100 commercial tape-outs using DSO.ai): news.synopsys.com/2023-02-07-... (<https://news.synopsys.com/2023-02-07-AI-designed-Chips-Reach-Scale-with-First-100-Commercial-Tape-outs-Using-Synopsys-Technology>)

16. [16] Cadence Cerebrus product page/case-study hub: [cadence.com/.../cerebrus-intelligent-chip-explorer.html](https://www.cadence.com/en_US/home/tools/digital-design-and-signoff/soc-implementation-and-floorplanning/cerebrus-intelligent-chip-explorer.html) (https://www.cadence.com/en_US/home/tools/digital-design-and-signoff/soc-implementation-and-floorplanning/cerebrus-intelligent-chip-explorer.html)
17. [17] Siemens Solido thought leadership on AI acceleration: blogs.sw.siemens.com/thought-leadership/2023/02/14/... (<https://blogs.sw.siemens.com/thought-leadership/2023/02/14/reshaping-the-ic-validation-and-characterization-industry-with-ai/>)
18. [18] The OpenROAD Project and repositories: github.com/The-OpenROAD-Project/OpenROAD (<https://github.com/The-OpenROAD-Project/OpenROAD>)
19. [19] CIRCT project site/repository: circt.llvm.org (<https://circt.llvm.org/>) and github.com/llvm/circt (<https://github.com/llvm/circt>)
20. [20] AlphaChip open-source implementation: github.com/google-research/circuit_training (https://github.com/google-research/circuit_training)
21. [21] Synopsys whitepaper on formal unreachability analysis: [synopsys.com/.../coverage-metrics-formal-unr-wp.html](https://www.synopsys.com/verification/resources/whitepapers/coverage-metrics-formal-unr-wp.html) (<https://www.synopsys.com/verification/resources/whitepapers/coverage-metrics-formal-unr-wp.html>)
22. [22] Huang et al., "Machine Learning for Electronic Design Automation: A Survey" (ACM TODAES, 2021): doi.org/10.1145/3447585 (<https://doi.org/10.1145/3447585>)
23. [23] GitHub Copilot supported models documentation: docs.github.com/copilot/reference/ai-models/supported-models (<https://docs.github.com/copilot/reference/ai-models/supported-models>)
24. [24] CVDP benchmark (2025), "Comprehensive Verilog Design Problems": arxiv.org/abs/2506.14074 (<https://arxiv.org/abs/2506.14074>)
25. [25] RTLLM benchmark and related open RTL evaluation efforts: arxiv.org/abs/2405.17378 (<https://arxiv.org/abs/2405.17378>)
26. [26] VeriCoder (2025), functional-correctness-guided RTL generation: arxiv.org/abs/2504.15659 (<https://arxiv.org/abs/2504.15659>)
27. [27] MCP4EDA (2025), tool-orchestrated RTL-to-GDSII automation research: arxiv.org/abs/2507.19570 (<https://arxiv.org/abs/2507.19570>)
28. [28] AutoEDA (2025), microservice-based LLM agents for EDA flow automation: arxiv.org/abs/2508.01012 (<https://arxiv.org/abs/2508.01012>)
29. [29] Anthropic, "Introducing the Model Context Protocol" (2024): [anthropic.com/news/model-context-protocol](https://www.anthropic.com/news/model-context-protocol) (<https://www.anthropic.com/news/model-context-protocol>)
30. [30] Survey: "The Dawn of Agentic EDA" (2025): arxiv.org/abs/2512.23189 (<https://arxiv.org/abs/2512.23189>)

Closing Thoughts

The semiconductor industry is in the early stages of a transformation as fundamental as the shift from schematic capture to RTL synthesis in the 1990s. That transition did not eliminate designers—it elevated them from drawing gates to describing behavior, enabling a leap in productivity that made billion-transistor SoCs possible. The AI transition will be similar: engineers will move from writing every line of RTL and every test stimulus to specifying intent, reviewing AI-generated implementations, and focusing their expertise on the architectural and verification decisions that machines cannot yet make.

The timeline for this transformation is measured in years, not months. Today's tools are powerful but narrow. LLMs generate useful RTL drafts but cannot verify them. RL agents close coverage gaps but require careful reward engineering. Commercial tools optimize PPA but within the parameter spaces their vendors have defined. The full vision—AI systems that understand a specification, generate the design, verify it, and optimize the physical implementation—remains a research challenge.

But the direction is clear, and the engineers and teams that begin building competency now will have a decisive advantage. The gap between "using AI" and "not using AI" in chip design will widen with each generation of models and tools. Start small. Pick one workflow. Collect data. Experiment. Iterate. The best time to begin was two years ago; the second-best time is today.

AI will not replace chip designers. But chip designers who master AI tools will design better chips, faster, with fewer bugs—and they will define the next era of semiconductor engineering.