

# Trusting the Machine: How Secure is LLM-Generated RTL Code?

Zahin Ibnat, Paul E. Calzada, Dipayan Saha, Hasan Al-Shaikh, Sujan Kumar Saha, Jingbo Zhou, Farimah Farahmandi, Mark Tehranipoor

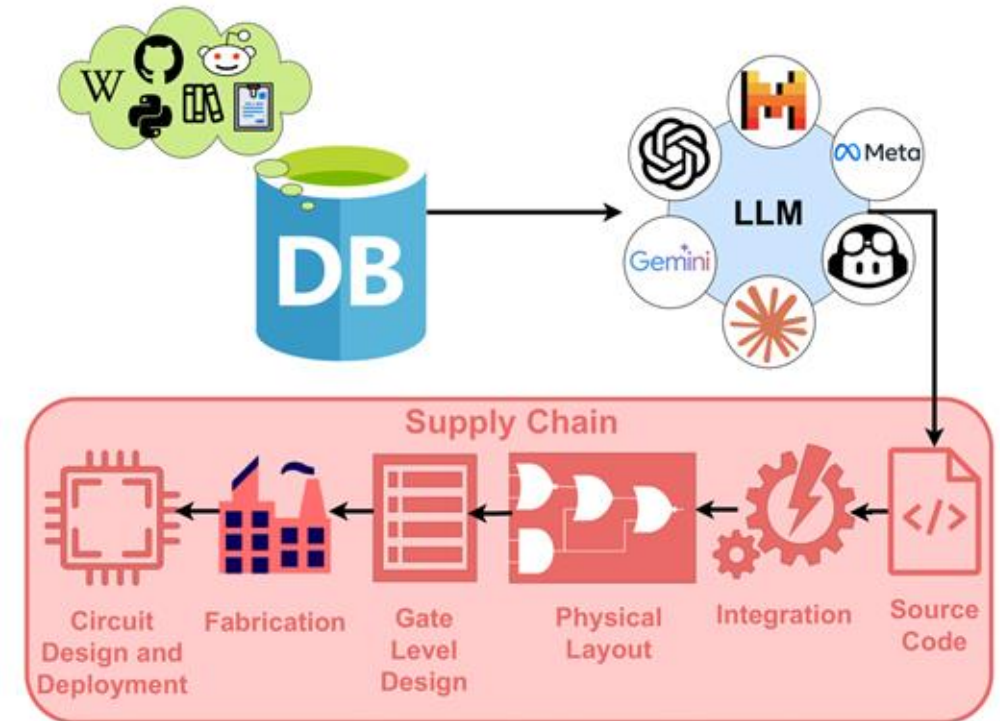
**Florida Institute for Cybersecurity Research  
Electrical and Computer Engineering Department**





Recent years, design complexities are growing despite shrinking time to market (TTM).

- Need for more custom designs with lower throughput lower energy.
- Aggressive TTM fosters need for automation to accelerate design flow.



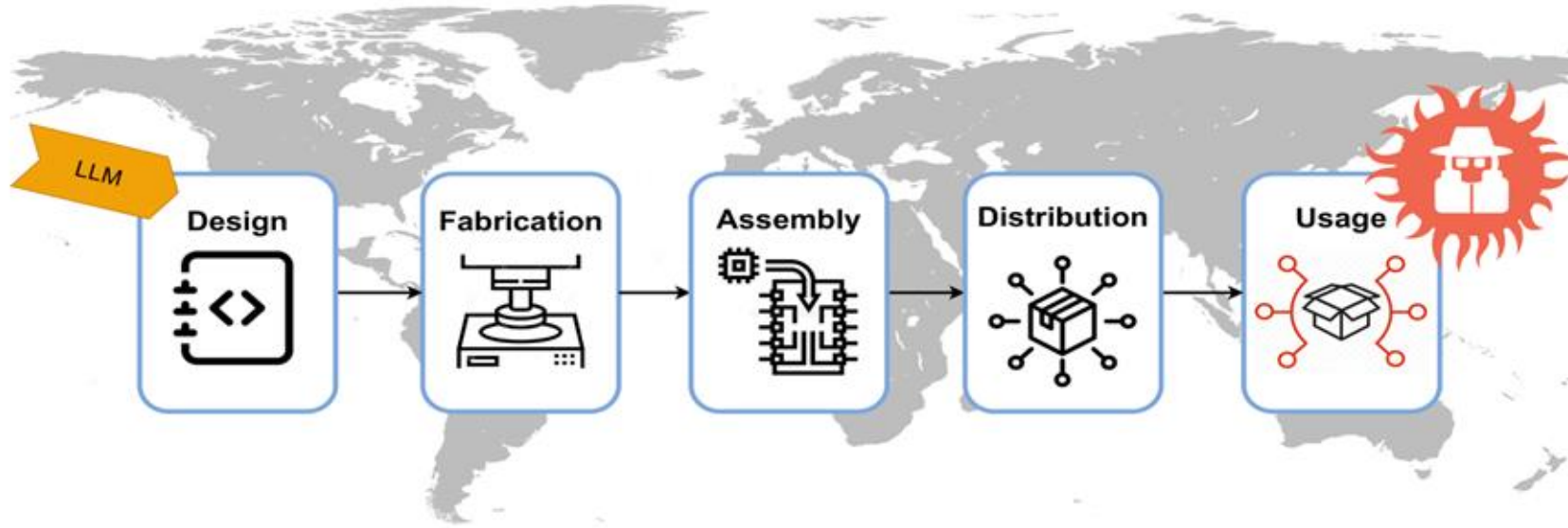
To automate segments of hardware design, companies are using LLMs.

- Ease design process through automation as designers prompt specs and receive LLM output as HDL.
- Increased industry use and research efforts for LLM generated RTL

# The Question of Security

With the desi  
could be the

- However,
- Are LLMs  
lack securi



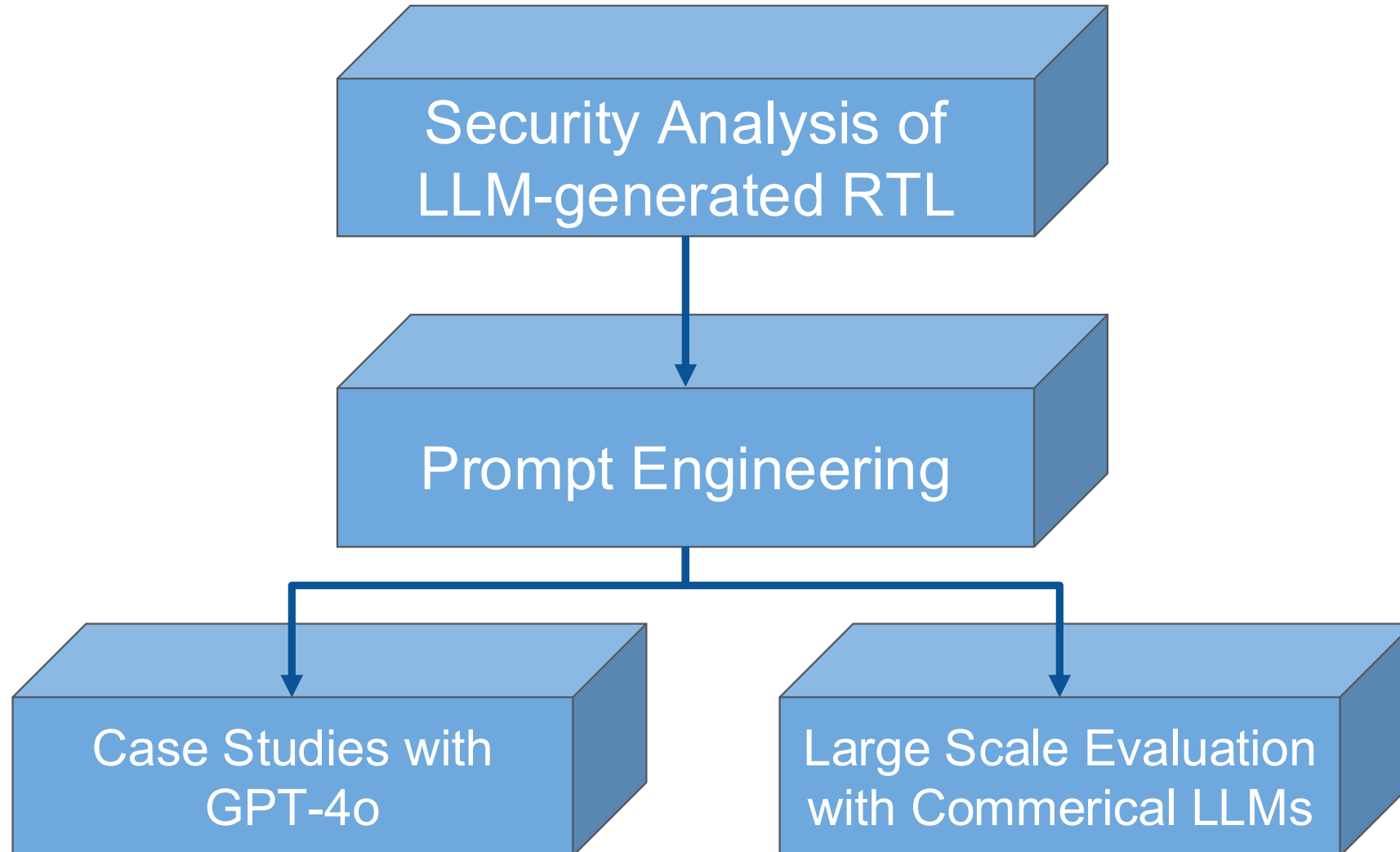
⇒ TTM, LLM

RTL code?  
which often

LLMs can generate designs with significant security flaws, such as improper handling of registers and missing access controls.

Of 42 LLM-generated designs, we found that 74% contained at least one security vulnerability.

**Threat Model:** An attacker could exploit these vulnerabilities to extract keys, disable modules, or trigger device failures after deployment.



Our **basic prompts** were designed to be straightforward and reflect how an engineer might use an LLM under aggressive TTM constraints.

You are an expert Verilog engineer specializing in digital design. You write clean, synthesizable, and minimal Verilog code. Write a basic 4 stage shift register in Verilog. It should:

- Have all registers start at 0
- Reset back 0 when 'rst' is asserted.
- Shift data sequentially through the registers every clock cycle

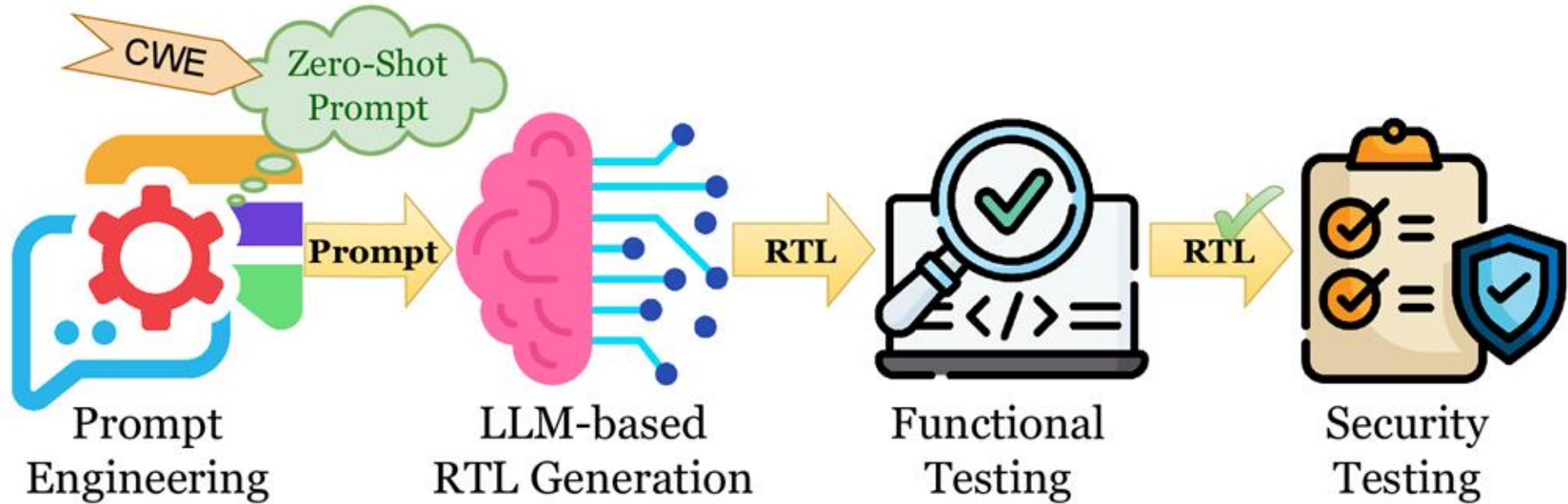
Our **complex prompts** build on the basic ones by adding a crucial layer: explicit security instructions based on CWEs, to assess the security understanding of LLMs.

You are an expert Verilog engineer specializing in digital design. You write clean, synthesizable, and minimal Verilog code. Write a basic 4 stage shift register in Verilog. It should:

- Have all registers start at 0
- Reset back 0 when 'rst' is asserted.
- Shift data sequentially through the registers every clock cycle
- Avoid CWE-1252: CPU Hardware Not Configured to Support Exclusivity of Write and Execute Operations
- Avoid CWE-190: Integer Overflow or Wraparound

This two-prompt approach allowed us to directly compare a purely functional request with a security-aware request and see how the LLM's output changed.

# Case Studies Setup and Flow



**Goal: Evaluate 8 different designs spanning design size and class for hardware vulnerabilities.**

Basic and Complex Prompting - Can prompting alone help LLMs generate secure RTL?

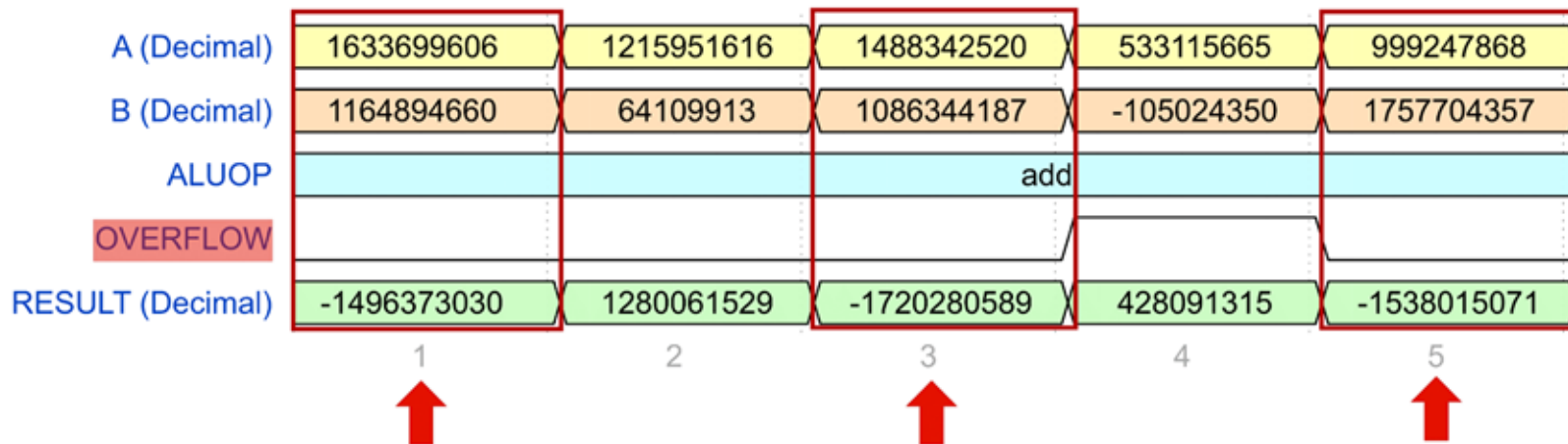
LLM-based RTL Generation - Can commercial models like GPT-4o generate secure code?

Functional Testing - Test vulnerabilities of RTL that would function in larger designs

Security Testing - *Expose Hardware Security Vulnerabilities in LLM-generated RTL*

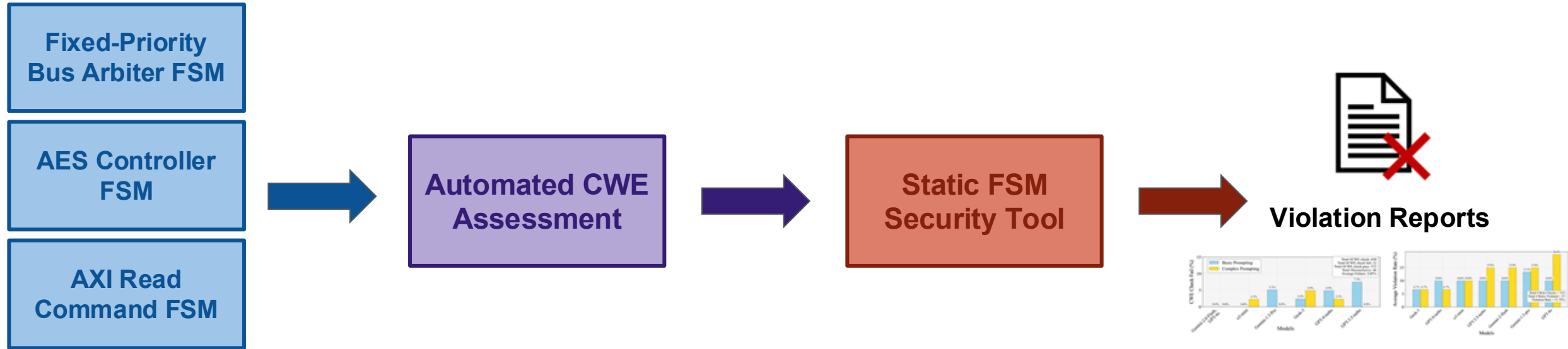
# Uncovering Specific Security Flaws: Case Studies

Design	Size	Description	Vulnerabilities	Prompted CWEs
RAM	Medium	RAM with write enable and reset	Simulation Failure: Lacked a read enable, causing the read bus to continuously output valid data and allowing unauthorized access.	CWE-226, CWE-1260
MIPS32 ALU	Large	ALU for 32-bit MIPS ISA CPU	Simulation Failure: The interface included only a single overflow signal, failing to detect signed arithmetic overflow as instructed.	CWE-190, CWE-1254
Counter	Small	8-bit counter with reset	Formal Verification Failure: The basic prompt failed the property assertion for proper overflow handling. When prompted with CWEs, the design's entire functionality changed to a saturating counter instead of fixing the overflow reporting	CWE-1221, CWE-190
FSM	Medium	8-state FSM with 3 state encoding bits	Formal Verification Failure: Unstable timing, switching states too early/operations done before their state	CWE-1221, CWE-1245



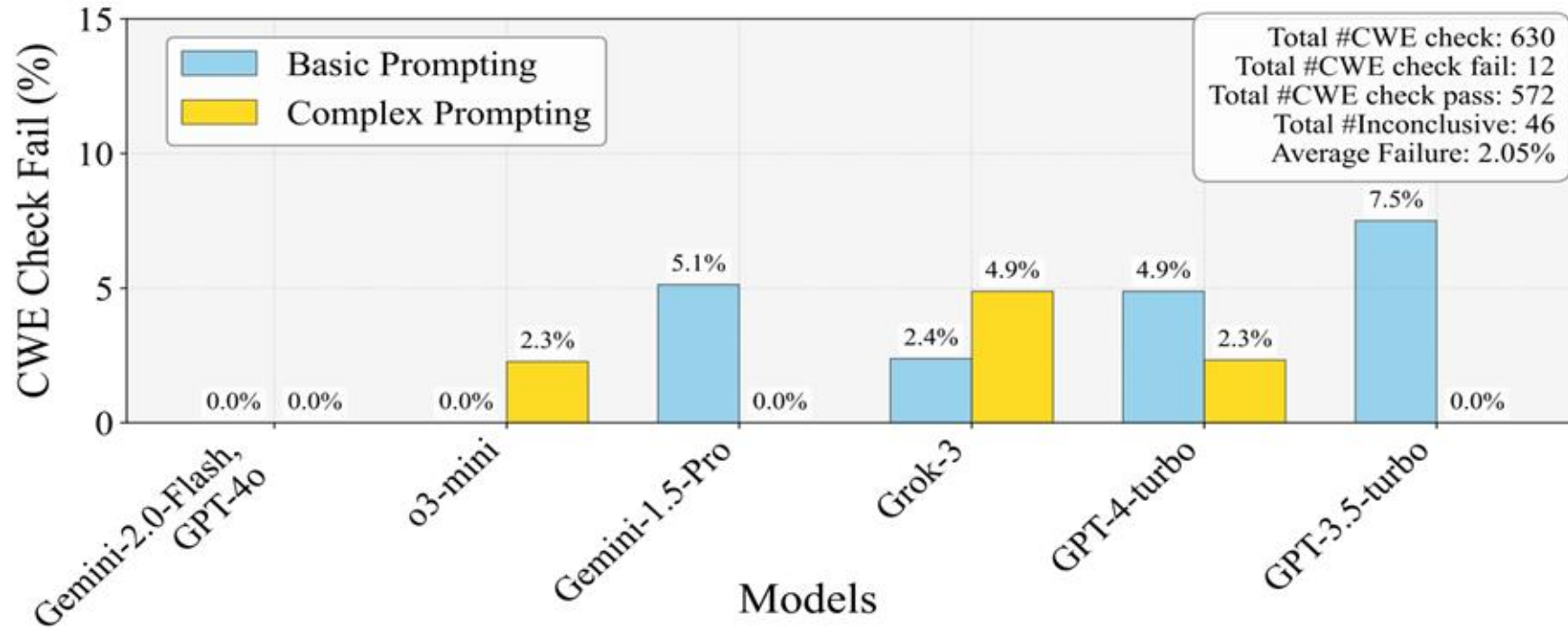
Overflow detection was explicitly prompted within the ALU, but signed arithmetic had no overflow flag raised.

# Large Scale Evaluation of LLM-Generated RTL



Performed both **basic and complex prompting** across 7 different commercial LLMs.  
Generated 3 widely used complex **FSM designs** -> 42 generated RTL codes  
**Automated CWE Assessment** for the designs against **15 hardware security CWEs**  
**Static FSM Security Tool** checked the designs against **9 specialized security rules**.

# CWE Violation Rate

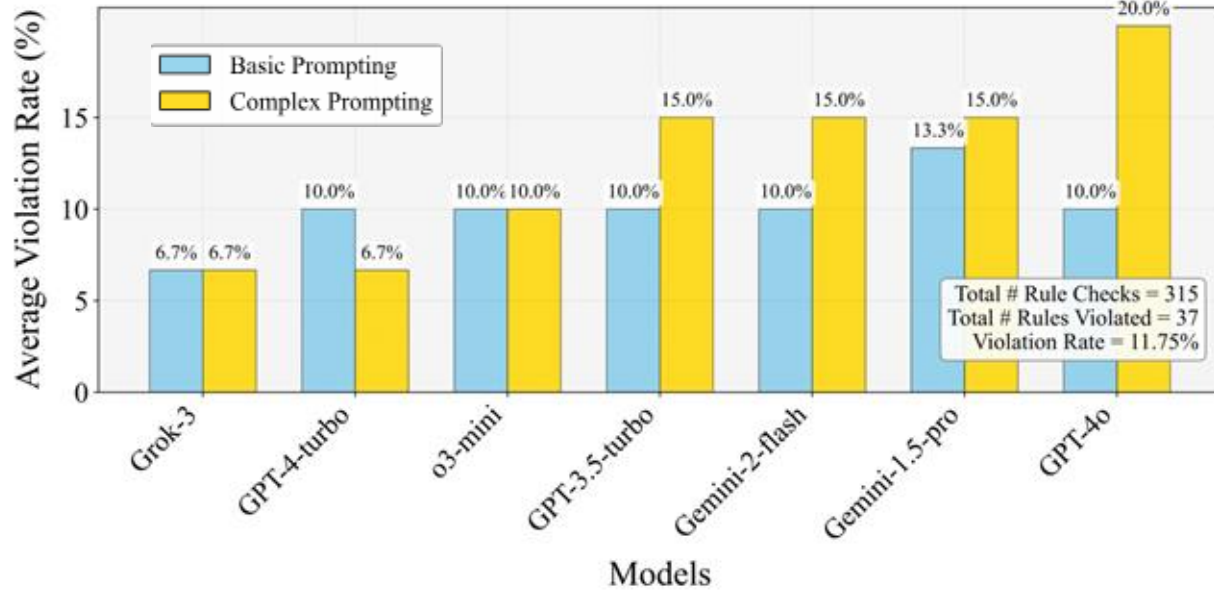


For general CWEs, complex prompting helped some models improve (largest impact on GPT-3.5-turbo), and newer models like GPT-4o and Gemini-2.0-Flash achieved 0% CWE failures.

Older or smaller models, which likely lack inherent knowledge of secure hardware design principles, can benefit from complex prompting.

**Overall average CWE failure rate across all designs and models remains relatively low at 2.05%, regardless of the prompting strategy applied.**

# FSM Rules Violation Rate



The top-performing models showed significantly higher violation rates on the FSM-specific security checks.

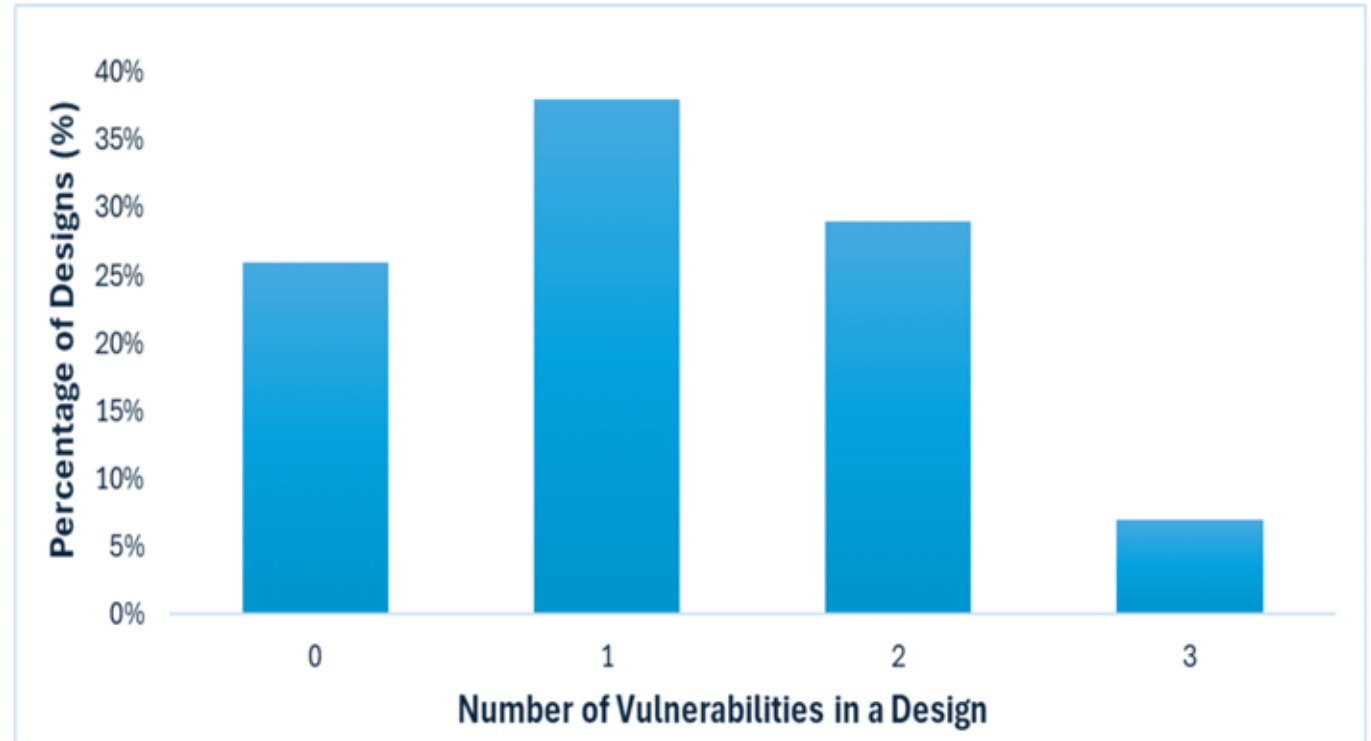
While LLMs may have some knowledge of widely documented CWEs, they lack a deeper understanding of FSM-level security best practices. Complex prompting can help reduce general vulnerabilities, but may introduce **structural design issues** due to increased design complexity or unnecessary transitions added during generation.

Rule	Description
R1	All unused states of a control FSM should be handled through the default statement in the RTL description.
R2	When state transition occurs between two consecutive unprotected states, Fault-Injection Feasibility (FIF) metric should be '0'.
R3	The reset state of a control FSM should be encoded to a particular defined value.
R4	Each state of a control FSM should be encoded uniquely.
R5	Control input signals of an FSM should be driven by trusted input ports only.
R6	Dead states should be absent in the extracted state transition graph (STG) of a control FSM from the RTL description.
R7	Unreachable states should not exist in the obtained state transition graph of a control FSM from the RTL description.
R8	States with static deadlock conditions should not be present in the extracted state transition graph of a control FSM.
R9	Groups of states with a dynamic deadlock loop that includes the protected state should be absent in the STG of a control FSM.

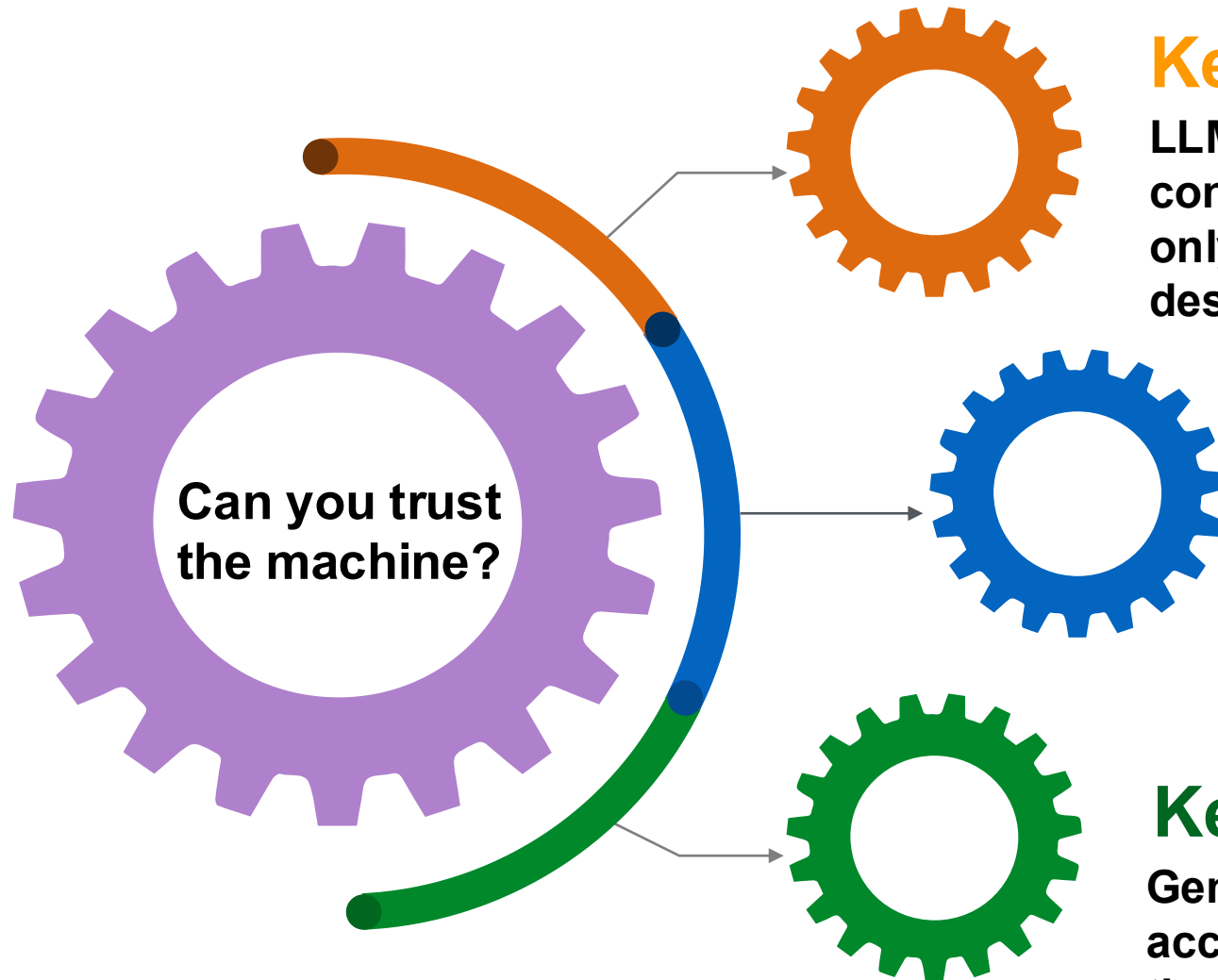
**Only 26% of the 42 generated designs were completely free of vulnerabilities.**

**The remaining 74% contained at least one security issue.**

**38% of designs had one vulnerability, 29% had two, and 7% had three.**



While LLMs can generate functionally correct designs, security correctness remains a challenge.



## Key Takeaway T1

LLMs can understand and apply security concerns from complex prompts, but typically only when generating simple or small-scale designs.

## Key Takeaway T2

Truly secure hardware design requires LLMs to have built-in security and programming understanding, rather than relying on external prompt guidance.

## Key Takeaway T3

Generated RTL from popular LLMs may fail to account for security-critical edge cases unless the design-specific interface and mechanisms are explicitly spelled out in the prompt.

While LLMs excel at automating IP design, they lack the inherent security awareness of a Verilog engineer. The models do not possess security reasoning, leading to IP vulnerabilities that can propagate risks across an entire System-on-Chip (SoC).

Although not all tested designs are security-critical in isolation, vulnerabilities in common RTL primitives can escalate risk when reused or integrated into security-sensitive SoCs.

Achieving secure RTL generation will require enhanced fine-tuning with verified datasets and other advanced AI-driven security approaches to ensure trustworthy hardware designs.

