

SENTAUR: Security EnhaNced Trojan Assessment Using LLMs Against Undesirable Revisions

Jitendra Bhandari, Rajat Sadhukhan, Prashanth Krishnamurthy, Farshad Khorrami, Ramesh Karri
New York University, New York, USA

Abstract—A globally distributed IC supply chain brings risks due to untrusted third parties. The risks span inadvertent use of hardware Trojan (HT), inserted Intellectual Property (3P-IP) or Electronic Design Automation (EDA) flows. HT can introduce stealthy HT behavior, prevent an IC work as intended, or leak sensitive data via side channels. To counter HTs, rapidly examining HT scenarios is a key requirement. While Trust-Hub benchmarks are a good starting point to assess defenses, they encompass a small subset of manually created HTs within the expanse of HT designs. Further, the HTs may disappear during synthesis. We propose a large language model (LLM) framework SENTAUR to generate a suite of legitimate HTs for a Register Transfer Level (RTL) design by learning its specifications, descriptions, and natural language descriptions of HT effects. Existing tools and benchmarks are limited; they need a learning period to construct an ML model to mimic the threat model and are difficult to reproduce. SENTAUR can swiftly produce HT instances by leveraging LLMs without any learning period and sanitizing the HTs facilitating their rapid assessment. Evaluation of SENTAUR involved generating effective, synthesizable, and practical HTs from TrustHub and elsewhere, investigating impacts of payloads/triggers at the RTL. While our evaluation focused on HT insertion, SENTAUR can generalize to automatically transform an RTL code to have defined functional modifications.

Index Terms—Hardware Trojan Detection, Hardware Security, Golden Reference-Free, Large Language Model

I. INTRODUCTION

Chip manufacturers face increasing challenges due to the scale and complexity of System-on-Chip (SoC) designs specifically targeted for modern embedded systems and Internet-of-Things (IoT) devices. Consequently, SoC designers, under time-to-market pressures and resource limitations, have turned to outsourcing hardware designs and utilizing Third-Party Electronic Design Automation tools or Intellectual Property cores (combined 3P-EDA/IP) from global vendors. This outsourcing presents numerous benefits, including cost reduction by minimizing internal team overheads, utilization of specialized skills through specialist firms for efficient resource allocation, faster development for crucial time-to-market scenarios, and access to a wider global talent pool with diverse skill sets.

However, the security and reliability of 3P-EDA/IPs remain uncertain, posing risks when relying on untrusted IPs and EDA tools. This raises concerns about potential malicious alterations within an IC, capable of causing side-channel leakages of sensitive data, functional changes, performance reduction, or denial of service (DoS). These malicious hardware alterations by unauthorized entities within the IC supply chain are commonly known as Hardware Trojans (HTs) [1],

which can even lead to the extraction of secret keys, enabling an adversary to modify the chip’s configurations and gain full control of the chip [2]. The standard design flow of an SoC commencing from design-level specification to physical IP development in a third-party environment is shown in Figure I (a). It highlights various stages within the IC supply chain identified as possible points where HTs could be injected. The adaptability of IP cores at higher abstraction levels in the design cycle makes it easy for attackers to insert HTs. Hence, developing detection and countermeasures against HTs at the design phase i.e. at the softIP level is critical [3], [4]; their elimination becomes costly in later stages. To measure the effectiveness of a countermeasure, designers need the ability to swiftly investigate the attack space susceptible to HT insertion for that design [5].

The HT benchmarks on Trust-Hub [6] are static and limited concerning the diverse nature of SoC development. According to [7] most Trust-Hub benchmarks rely on unrealistic assumptions. They show that 3-out-of-83 HT benchmarks are effective. Further, FPGA/ASIC development process entails designing control and data paths requiring deep hardware expertise. When modifications/additions are frequent to accelerate an algorithm, it can result in frequent obsolescence problems. Unlike software that can be updated or recompiled, IC designs cannot incorporate changes seamlessly. There is a critical need for a platform-independent framework to add functions to an IC design at RTL.

A. Key Contributions

SENTAUR is an LLM-based HT assessment flow shown in Figure I (b). Given a design specification and corresponding RTL, we query SENTAUR to generate HTs based on a HT trigger class, and payload class. Our contributions are:

- A novel use of conversational LLMs to generate RTL that is synthesizable using a HT netlist from Trust-Hub. A recent study [7] shows that most HTs are not correct and effective after synthesis. We use LLMs to insert HTs which after synthesis persist.
- SENTAUR is a flexible, versatile, and platform-independent LLM-based toolchain for inserting design templates given functional descriptions of the design. It can be used by an adversary to insert or analyze the effect of HTs and by a designer to plug in templates.
- Validate SENTAUR flow from an attacker view using the Xilinx platform. We inserted different HTs in the RTL and validated the designs using the Xilinx FPGA.

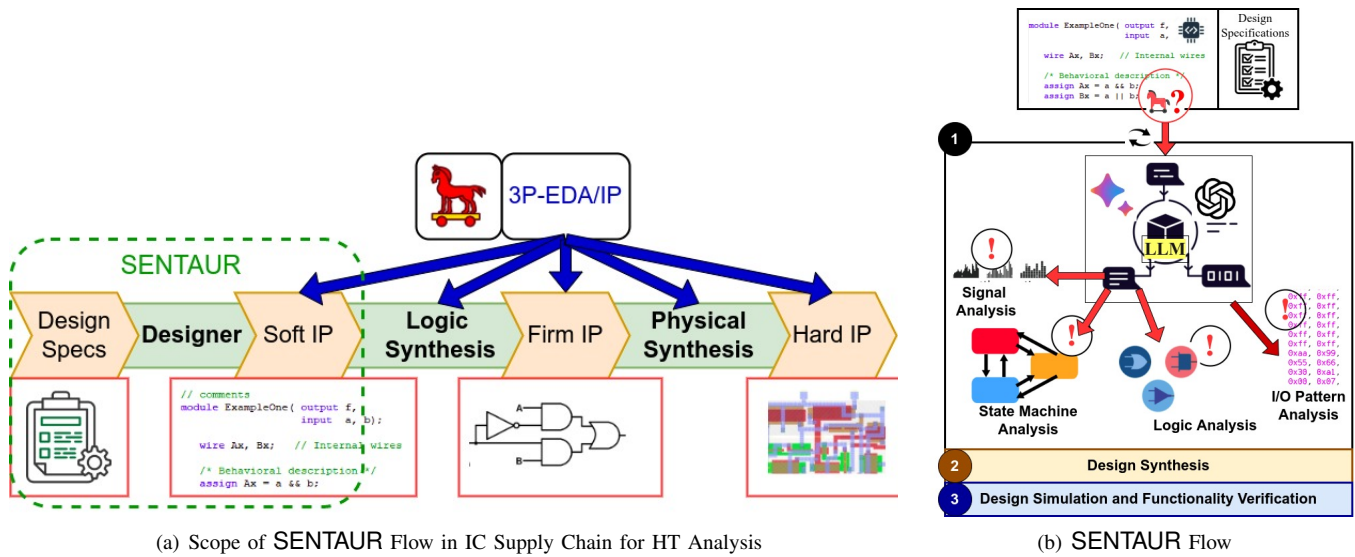


Fig. 1. SENTAUR Flow and its scope in IC Supply Chain.

We proposed a mechanism to sanitize and generate HTs from Trust-Hub through SENTAUR such that functionally correct HTs post-synthesis are generated.

B. Related Works

In this section, we will discuss the state-of-art research in HT insertion and detection. The very first work in this direction is the development of an extensive repository of HTs available at Trust-Hub [6]. However, Trust-Hub is restricted to the number of HT circuits that cover only a fraction of the potential landscape for inserting HTs in digital circuits, thereby restricting the development of varied countermeasures. To overcome this limitation TAIN [8] tool is proposed where HT insertion is done at the various stages of the design cycle. However, the tool anticipates that the user will choose the trigger nets based on recommendations provided by the tool itself. In [9] the authors proposed an automated tool TRIT to insert HTs in a design by configuring various parameters such as the number of trigger nets, the count of rare nets among these triggers, rare-net threshold determined through signal probability of nets, and the selection of payload. Despite expanding the range of inserted HTs, the TRIT methodology cannot identify the best trigger and payload nets. The work [10] proposes a flow that explores Trojans in physical design layouts restricted to ASIC layouts and applicable at the backend stage. Similarly, in [11] Trojan space is explored at the backend stages of FPGA design flow. In [12], the authors proposed a flow called HAL that inserts countermeasures against Trojan attacks but doesn't explore the Trojan space.

Yu et al. [13] proposed a methodology that identifies rare nets using the transition probability of nets to insert HTs. MIMIC [14] leverages ML to insert HTs. MIMIC ML model was developed by extracting 16 functional and structural attributes from existing HT samples, creating numerous HTs tailored for specific designs. MIMIC process is intricate,

TABLE I
COMPARISON OF PROPOSED SENTAUR WITH STATE-OF-ART HT INSERTION TOOLS

Work	Effort	Trigger Net Selection	Learning Time Required
TAIN [8]	Manual	User Selects Trigger Nets	NA
TRIT [9]	Automated	Signal Probability of Nets	NA
Yu et al. [13]	Automated	Transition Probability of Nets	NA
MIMIC [14]	Automated	supervised and generative ML to extract features of Nets	Yes
HT Playground [15]	Automated	RL based location exploration	Yes
SENTAUR (This Work)	Automated	User Prompts	No

involves multiple stages, and requires extensive learning time to train the model. Trojan Playground [15] functions the same way where Reinforcement Learning (RL) is used requiring extensive training of the model. We propose an LLM-based tool flow SENTAUR that does not involve a time-consuming training process and is user-friendly in generating extensive set of HTs and insert them into an RTL design. A summary of SENTAUR tool and its benefits are shown in Table I.

The rest of the paper is organized as follows: Section II-A summarizes the HT threat model, II-B motivates this study with a use case, followed by Section II-C and Section II-D, which details the SENTAUR flow and capabilities respectively. Section III expounds on the experimental configuration and outcomes, and lastly, Section IV presents the conclusions and potential future directions.

II. SENTAUR: TOOL FLOW AND CAPABILITIES

A. Hardware Trojan Threat Model

Our proposed flow SENTAUR aims to assess malicious possible HT circuits in RTL code, concerning trigger, and payload design. SENTAUR is capable of analyzing or identifying trigger-based or continuously active HTs in an RTL code given its specification, including those with a payload circuit intended to alter functionality, reduce performance, leak sensitive data, or disrupt service. Our approach focuses on HT insertion during the design stage through scenarios involving deliberate manual manipulation by untrustworthy 3PIP vendors or 3P-EDA tools targeting the RTL netlist.

B. Motivating Use Case

Krieg et al. [7] argue that the assumptions underlying TrustHUB HT benchmarks are unlikely to apply in real-world scenarios. In this section, we will show LLMs can help generate a synthesizable code that reproduces the result as intended. For this, we analyzed the AES-T800 benchmark for our experiment which is also used in [7] to demonstrate that the benchmark violates the correctness and persistence property after synthesizing the netlist. HT insertion in the AES-T800 benchmark detects a predetermined sequence of input plaintext and then illicitly transmits the AES secret key via a hidden communication channel. However, in [7] it has been found that post-synthesis the trigger signal rises, but it appears that it doesn't occur due to the predefined sequence. Instead, it rises right after the circuit reset goes off. Hence using these benchmarks requires manual intervention by hardware experts. They sanitize the designs to make them synthesizable ensuring usability, correctness, and consistency.

We leveraged the power of LLM to generate a synthesizable code efficiently. We chose GPT-4 in our experiment. On prompting GPT-4 as shown below the RTL is generated as shown in Listing 1. We synthesized the code using Xilinx Vivado and simulated the post-synthesized using the benchmark from Trust-Hub. The simulated result of the post-synthesized netlist is shown in Figure 2; after detecting the sequence of four input states the trigger signal goes high.

Prompt: ..., once a predefined sequence of 4 input plaintext (state) is observed, alert by setting a signal (Tj_Trig) Write a hardware module in Verilog that is synthesizable.

LLM Generated Synthesizable Code:

```
module sequence_detector(input wire clk, input wire rst,
    input wire [127:0] state, output reg Tj_Trig);

localparam [2:0] IDLE = 3'b000, SEQ1 = 3'b001, SEQ2 = 3'
b010, SEQ3 = 3'b011, SEQ4 = 3'b100, TRIGGER =
3'b101; // Define states
reg [2:0] current_state, next_state;
always @(*) begin // Comb. logic for state transitions
    case (current_state)
        IDLE:
            if (state == 128'
h3243f6a8_885a308d_313198a2_e0370734)
                next_state = SEQ1;
            else
                next_state = IDLE;
```

```
SEQ1:
    if (state == 128'
h00112233_44556677_8899aabb_ccddeeff)
        next_state = SEQ2;
    else
        next_state = IDLE;
    ...

SEQ4:
    next_state = TRIGGER;

TRIGGER:
    next_state = IDLE;
// Default case provided by LLM to not infer latch post-
synthesis
default:
    next_state = IDLE;
endcase
end

always @(posedge clk) begin
    if (rst) // Sequential logic for state updating
        current_state <= IDLE;
    else
        current_state <= next_state;
end

// Output logic
always @(current_state) begin
    Tj_Trig = (current_state == TRIGGER) ? 1'b1 : 1'b0;
end
endmodule
```

Listing 1. Synthesizable AEST-800 HT Benchmark adhering to Correctness and Persistence

We further expanded the scope and result of this experiment to develop an LLM-based flow SENTAUR that creates a diverse set of valid HTs based on the functionality of a particular design at the RTL level adhering to correctness post-synthesis. Our tool can be used as a more generic version where it can offer a versatile framework that facilitates automated alterations of provided RTL code to implement specific and predefined functional modifications.

C. Tool Flow

In the last section, we have seen how LLM could generate an RTL code of the Trust-Hub [6] Trojan benchmark that after synthesis gives functionally correct and intended results. We verified the result using gate-level simulation against the Trojan benchmark provided in Trust-Hub. In this section, we will formalize the flow and introduce an automation framework SENTAUR that will generate RTL codes given the specification of the RTL design through an LLM as shown in Figure I (b). Our proposed flow is a three-stage process

- 1) *RTL Generation/Analysis using Specification:* The function description of the design under consideration is given to the LLM to generate variations in functionality concerning logic, state machine, I/O pattern, and signal analyses. These functionalities are generated based on some trigger condition specified to LLM based upon time, input-output patterns, or physical conditions. Given that under the purview of HT insertion the generated functionality varies under different trigger conditions, the final effect can be seen as either DoS, leakage of signal values, or degradation of design performance.

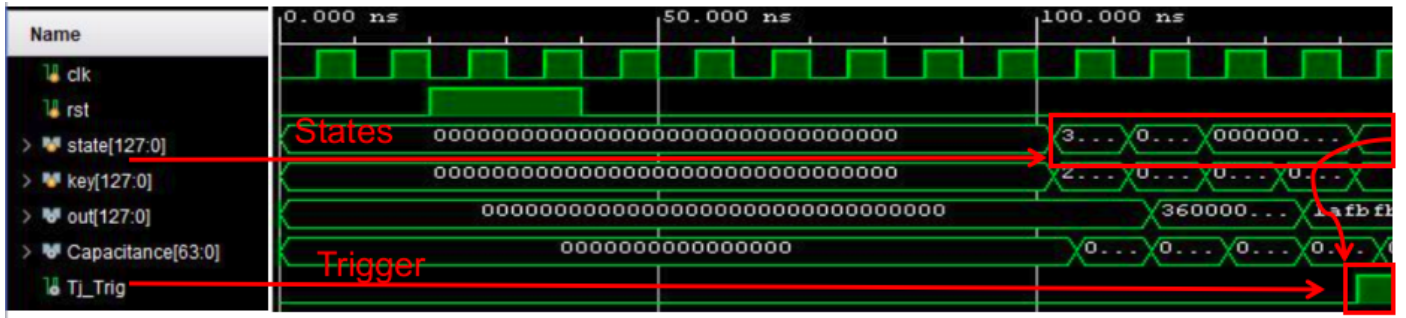


Fig. 2. Post-synthesis Simulations of AES-T800

From a designer’s perspective, this feature can be used to plug in templates in a design, and from a defender’s perspective, one can use this tool to insert countermeasures to protect against physical attacks. The LLM is directed to generate behavioral descriptions that are synthesizable ensuring functional correctness and persistency at the later stages of the IC design process. Modifications to the behavioral model also enable the tool to be independent of ASIC and FPGA platforms making it a universal tool. The LLM-based generation also enables the flow to be agnostic to vendor-specific file formats, reducing the need for programming or hardware expertise.

- 2) *Design Synthesis*: In this stage the generated RTL design from stage 1 is synthesized using commercial or open-source synthesis tools to generate a gate-level netlist.
- 3) *Design Verification*: The synthesized netlist is verified using a commercial or open-source simulation tool to ensure the desired functionality and correctness.

D. SENTAUR Tool Capabilities

We will delve into the details of two capabilities in this section and examine how the tool can be utilized in scenarios involving attackers, defenders, and designers.

1) *Signal Declarations and Connections*: SENTAUR possesses the capability to incorporate, modify, or eliminate signals, enabling the introduction of new signals tailored for HT-related functions or the alteration of existing signals to include malicious logic. The tool’s manipulation of signal connections can be utilized to establish covert communication pathways between the inserted module and external entities. The tool supports various signal modifications:

Input/Output Signals: The tool augments module declaration by incorporating input/output signals based on user-defined parameters. It caters to scenarios where signals are omitted.

Join Signals: These signals unify signals. By replacing the original signal name with the joint signal name, the tool facilitates adding these combined signals into the design. Route and join signals are convenient pathways for transmitting signals that trigger HTs or instructions activating malicious actions.

Add-on Signals: specified in the extra signals parameter integrate into the module declaration.

2) *File Manipulation*: SENTAUR reads the RTL netlist, applies modifications, and saves the modified RTL netlist to the

file. This aligns with an attacker’s objective of surreptitiously injecting HT logic without raising suspicion.

SENTAUR offers valuable functionality for modifying HDL files in designs. It can rename modules, manage signal declarations and connections, and manipulate RTL netlists. Designers can customize FPGA/ASIC designs to meet requirements, integrate subsystems, and streamline design processes by utilizing SENTAUR. Its features grant attackers the ability to obfuscate and embed malicious HTs within FPGA/ASIC designs. Additionally, SENTAUR aids assessing vulnerabilities in FPGA/ASIC designs and studying impacts of HTs.

III. EXPERIMENTAL RESULTS

A. Setup

For our study, we integrated GPT-4 to our proposed flow SENTAUR in generating and assessing various HT functionalities. Our RTL benchmark set consists of HTs from Trust-Hub viz. AES-T600, AES-T800, AES-T900, and standard IP dual-port RAM from Xilinx integrated with Processing System (PS) along with additional logic to control the IP. The motivation to choose dual-port RAM for our experiments is that it finds applications in various fields, including digital signal processing, networking, multiprocessing systems, and high-performance computing. Dual-port RAM refers to a type of memory structure that enables two separate read and write operations to occur simultaneously requiring two different clocks. This memory setup allows for independent and simultaneous access to the data stored within it. It’s commonly utilized in scenarios where two different processes or modules need access to the same memory resource without causing conflicts or delays. For testing the changes done on RTL by our proposed flow SENTAUR, we employed a real-world Xilinx FPGA of the Zynq-7000 family with device part xc7z020clg400-1. We built a Linux image for the PS with *jupyter notebook* interface compatibility, to control and program the functionality of our SoC design using Python code, with the capability to read and write the memory. Figure 3 and Figure 4 shows the design implemented in the FPGA without and with the HT presence, respectively. This validates that the modification done on top of the RTL by the GPT-4 on the RTL is synthesizable.

TABLE II
 RESULT ON THE OVERHEAD ASSOCIATED WITH {USING GPT-4 AND MANUAL} AND THE ASSESSMENT USING SENTAUR FRAMEWORK.
 NOTE: ● SIGNIFIES RAISED CONCERN BY THE LLM AND ○ SIGNIFIES NOT FLAGGED BY THE LLM, RESPECTIVELY.

Trigger	Description	Effect	GPT-4		Manual		Assessment			
			LUT	FF	LUT	FF	I/O	FSM	Logic	Signal
Time-based	When the count reaches in between 50 and 200.	DOS	791	1261	777	1080	○	○	●	●
		Perf. Degrade.	846	1261	829	1080	○	○	●	●
		Inf. Leak	798	1261	779	1080	○	○	●	●
Logic-based	When the value of the data is between 1000 and 2000.	DOS	772	1072	774	1072	○	○	●	●
		Inf. Leak	778	1072	778	1072	○	○	●	●
Address	When the input address is in the range of 2000 and 3000.	DOS	776	1072	774	1072	○	○	●	●
		Perf. Degrade.	829	1072	825	1072	○	○	●	●
		Inf. Leak	781	1072	779	1072	○	○	●	●
State based	When the FSM detects the sequence 0x55, 0xAA and 0xFF.	DOS	785	1258	821	1076	○	●	●	●
		Perf. Degrade.	837	1258	853	1076	○	●	●	●
		Inf. Leak	793	1258	834	1076	○	●	●	●
Input based	When the # of writes done reaches 1000.	DOS	782	1104	776	1084	●	○	●	●
		Perf. Degrade.	834	1104	823	1084	●	○	●	●
		Inf. Leak	791	1104	781	1084	●	○	●	●
AES-T600 [6]	Detects a specific input plaintext	Inf. Leak	4588	4475	4682	4521	○	○	●	●
AES-T800 [6]	Checks for a predefined sequence of input plaintext is observed	Inf. Leak	4591	4483	4726	4525	○	●	●	●
AES-T900 [6]	After each 2^{128} encryptions, the gets activated	Inf. Leak	4771	4502	4834	4523	●	○	●	●

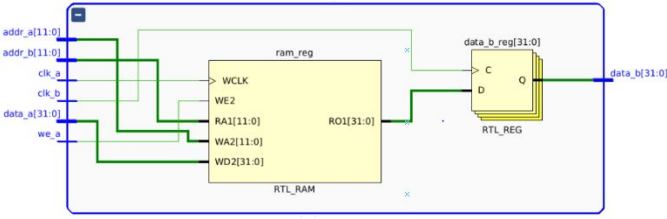


Fig. 3. Schematic of Dual-port RAM by Xilinx without HT.

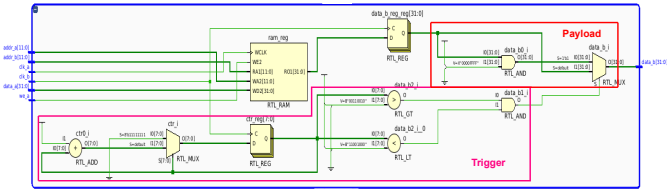


Fig. 4. Dual-port RAM with HT inserted (State-based) by SENTAUR.

B. Result I: Analysis on Generation of HTs

For this study, we took various combinations of triggers and effects for our analysis. Specifically, we took 5 triggers {time-based, logic-based, address, state-based, input-based} and 3 effects {Denial-of-service (DoS), Performance Degradation (Perf. Degrade), Information Leakage (Inf. Leak)}, respectively. Table II gives the overall taxonomy as well as the description of each trigger. For our effects:

- *Perf. Degrade* by introducing a dead band which compromises the performance in terms of output being 0 for some small interval at a regular instant (similar to having some delay in between 2 transmissions) with a slight impact in the amplitude as shown in Figure 5.
- *Info. Leak* signifies side-information leakage as shown in Figure 6 through the use of different channels than the one used for signal transmission, thus making an

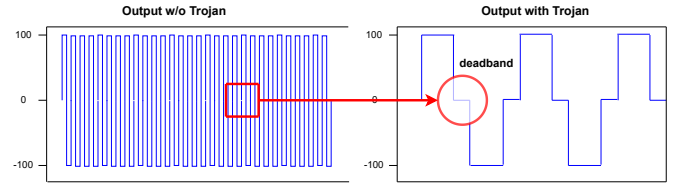


Fig. 5. Expected outputs with and without presence (Performance degradation by introducing deadband) when the trigger condition is met.

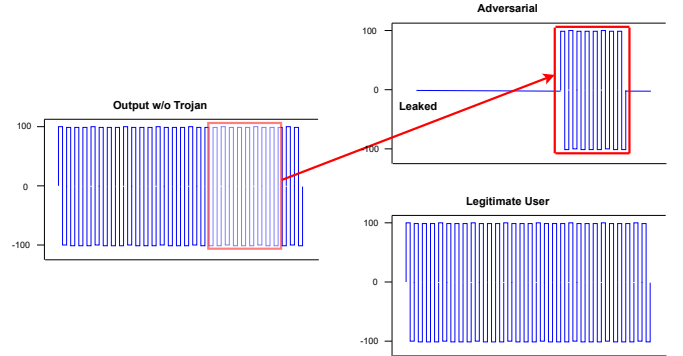


Fig. 6. Expected outputs with and without (Information leakage by copying the data to Adversarial) when the trigger condition is met.

adversary aware of the data being transmitted to some other users. This can compromise the privacy of a user.

- *DoS* by making the output 0 thus denying any service by not transmitting the required data as shown in Figure 7.

We generated the required changes in the RTL by appropriately prompting the GPT-4 with the information of the desired result. To compare the ones generated by the RTL and how a designer would have written them, we compared the LUTs and FFs required by each of them. Table II summarizes our result of the different overheads required in both scenarios demonstrating a comparable result when LLM is used. In addition, we took a few examples from the Trust-Hub [6]

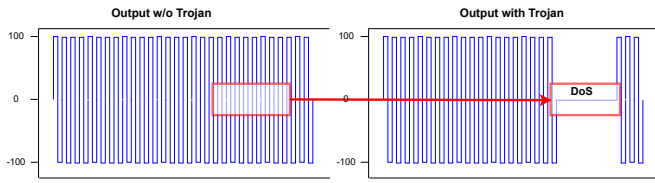


Fig. 7. Expected outputs with and without presence (Denial of Service by making the output 0) when the trigger condition is met.

specifically, AES-T{600,800,900} and did a similar analysis, where we compared the overhead associated while designed by a human and LLM, respectively.

C. Result II: Assessment of HTs

For the second set of analyses, we focused on manually written HT codes, using LLM (GPT-4), to evaluate these codes for malicious components or indicators of suspicious activity. To ensure a comprehensive assessment, we employed a variety of code combinations as detailed in Table II including examples from Trust-Hub [6]. A key strength of GPT-4 is its ability to summarize and analyze code, allowing it to identify potential vulnerabilities within the code. Crucial to this process is the use of carefully crafted prompts, which guide the LLMs in their analysis. We structured the prompts around four specific areas: logic, state machines, I/O pins, and signal analysis, aligning with the methodology depicted in Figure I(b). This targeted approach enabled the harnessing of GPT-4 analytical capabilities to detect and flag areas in the HT codes that might pose security risks.

- *I/O pin*: In this analysis, it looks for any I/O pins that are used for some sort of condition (or trigger) that can modify the result in the code.
- *State Machine*: In this, the analysis seeks to flag FSMs that look for a particular sequence and depending on that satisfy some condition that can affect the result. There can be many FSMs in a real-world RTL code, which from a human point of view become quite difficult to keep track of, so if LLM can point to only those FSMs where there is some suspicious behavior, it will be of tremendous help from a designer point of view.
- *Logic*: This is the most common flag detected by LLMs in various scenarios, indicating distinct logic in the code or dependency on specific conditions to activate. While this might trigger false alarms, focusing on thorough code scrutiny, especially concerning security, outweighs the potential risks of overlooking critical issues.
- *Signal*: For this, we look for any potentially vulnerable signal that can trigger some part of the inactive logic in normal operations. LLMs can do a great job in reporting only those signals where a human has to go through the whole code and follow the transition of the signal which becomes cumbersome in a large codebase.

Table II summarizes our findings under the column name ‘Assessment’. It shows the part of the code being flagged as potentially vulnerable by the LLM.

IV. CONCLUSION AND DISCUSSION

SENTAUR is a framework leveraging an LLM to generate a diverse set of legitimate HTs at the RTL. Unlike existing tools that require a learning period to replicate threat models, SENTAUR rapidly generates HT instances using LLMs, bypassing the learning phase. It effectively assesses Trust-Hub HTs, conducting comprehensive evaluations with practical use cases and Trust-Hub benchmarks, and exploring diverse impacts and trigger mechanisms at the RTL level. While our primary focus is HT insertion evaluation, SENTAUR also offers a flexible framework for automated modifications of RTL to integrate specific functionalities. Future direction involves extending this work to insert functionalities post-synthesized netlist.

REFERENCES

- [1] M. Tehranipoor and F. Koushanfar, “A survey of hardware trojan taxonomy and detection,” *IEEE Design & Test of Computers*, vol. 27, no. 1, pp. 10–25, 2010.
- [2] S. Skorobogatov and C. Woods, “Breakthrough silicon scanning discovers backdoor in military chip,” in *Cryptographic Hardware and Embedded Systems – CHES 2012*, E. Prouff and P. Schaumont, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 23–40.
- [3] R. Karri, J. Rajendran, K. Rosenfeld, and M. Tehranipoor, “Trustworthy hardware: Identifying and classifying hardware trojans,” *Computer*, vol. 43, no. 10, pp. 39–46, 2010.
- [4] V. R. Surabhi, P. Krishnamurthy, H. Amrouch, J. Henkel, R. Karri, and F. Khorrami, “Exposing hardware trojans in embedded platforms via short-term aging,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 11, pp. 3519–3530, 2020.
- [5] R. S. Chakraborty, S. Narasimhan, and S. Bhunia, “Hardware trojan: Threats and emerging solutions,” in *2009 IEEE International High Level Design Validation and Test Workshop*, 2009, pp. 166–171.
- [6] H. Salmani, M. Tehranipoor, and R. Karri, “On design vulnerability analysis and trust benchmarks development,” in *2013 IEEE 31st International Conference on Computer Design (ICCD)*, 2013, pp. 471–474.
- [7] C. Krieg, “Reflections on trusting trusthub,” in *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*, 2023, pp. 1–9.
- [8] V. Jyothi, P. Krishnamurthy, F. Khorrami, and R. Karri, “Taint: Tool for automated insertion of trojans,” in *2017 IEEE International Conference on Computer Design (ICCD)*, 2017, pp. 545–548.
- [9] J. Cruz, Y. Huang, P. Mishra, and S. Bhunia, “An automated configurable trojan insertion framework for dynamic trust benchmarks,” in *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2018, pp. 1598–1603.
- [10] J. Bhandari, J. Gopinath, M. Ashraf, J. Knechtel, and R. Karri, “Defending integrated circuit layouts,” *Cryptology ePrint Archive*, Paper 2023/205, 2023, <https://eprint.iacr.org/2023/205>. [Online]. Available: <https://eprint.iacr.org/2023/205>
- [11] J. Cruz, C. Posada, N. V. R. Masna, P. Chakraborty, P. Gaikwad, and S. Bhunia, “A framework for automated exploration of trojan attack space in fpga netlists,” *IEEE Transactions on Computers*, vol. 72, no. 10, pp. 2740–2751, 2023.
- [12] M. Fyrbiak, S. Wallat, P. Swierczynski, M. Hoffmann, S. Hoppach, M. Wilhelm, T. Weidlich, R. Tessier, and C. Paar, “Hal—the missing piece of the puzzle for hardware reverse engineering, trojan detection and insertion,” *IEEE Transactions on Dependable and Secure Computing*, vol. 16, no. 3, pp. 498–510, 2019.
- [13] A. Bhattacharyay, S. Yang, J. Cruz, P. Chakraborty, S. Bhunia, and T. Hoque, “An automated framework for board-level trojan benchmarking,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 42, no. 2, pp. 397–410, 2023.
- [14] J. Cruz, P. Gaikwad, A. Nair, P. Chakraborty, and S. Bhunia, “Automatic hardware trojan insertion using machine learning,” *arXiv preprint arXiv:2204.08580*, 2022.
- [15] A. Sarihi, A. Patooghy, P. Jamieson, and A.-H. A. Badawy, “Trojan playground: A reinforcement learning framework for hardware trojan insertion and detection,” *arXiv preprint arXiv:2305.09592*, 2023.