

This paper has been accepted at IEEE International Symposium on Hardware Oriented Security and Trust (HOST), 2025.

# SoCureLLM: An LLM-driven Approach for Large-Scale System-on-Chip Security Verification and Policy Generation

Shams Tarek, Dipayan Saha, Sujan Kumar Saha, Mark Tehranipoor, Farimah Farahmandi  
*Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL, USA*  
{shams.tarek, dsaha, sujansaha}@ufl.edu, {tehranipoor, farimah}@ece.ufl.edu

**Abstract**—Contemporary methods for hardware security verification struggle with adaptability, scalability, and availability due to the increasing complexity of the modern system-on-chips (SoCs). Large language models (LLMs) have emerged as a viable approach to address these shortcomings in security verification because of their natural language understanding, advanced reasoning, and knowledge transfer capabilities. However, their application to large designs is limited by inherent token limitation and memorization constraints. In this paper, we introduce SoCureLLM, an LLM-based framework that excels in identifying security vulnerabilities within SoC designs and creating a comprehensive security policy database. Our framework is adaptable and adept at processing varied, large-scale designs, overcoming the abovementioned issues of LLM. In evaluations, SoCureLLM detected 76.47% of security bugs across three vulnerable RISC-V SoCs, outperforming the state-of-the-art security verification methods. Furthermore, assessing three additional large-scale RISC-V SoC designs against various threat models led to the formulation of 84 novel security policies, enriching the security policy database. Previously requiring extensive manual effort to craft, these newly generated security policies can be used as guidelines for developing secured SoC designs.

**Index Terms**—Large Language Model, Hardware Security, Verification, Security Bug Detection, Policy Generation

## I. INTRODUCTION

Modern system-on-chips (SoCs) integrate third-party intellectual property (3PIP) from global vendors, making security verification crucial due to their complexity and the horizontal development process. Hardware security vulnerabilities may occur due to unintentional design mistakes, deliberate actions by malicious insiders, security-unaware optimization by CAD tools, or flaws within test and debug infrastructures. The rise in microarchitectural attacks highlights the need for robust security verification to protect assets such as sensitive user data, cryptographic keys, and configuration information. Hardware vulnerabilities can cause massive financial losses and product recalls in the semiconductor industry [1]. To address this, there is an industry shift towards a security development lifecycle (SDL) for SoCs [2]. Nonetheless, current verification tools

focus more on functional than security aspects, often missing critical vulnerabilities in the pre-silicon phase.

Hardware security vulnerabilities can be of different types, such as information leakage, access control violation, unauthorized memory access, fault injection attacks, and side-channel attacks. To address these vulnerabilities, recent research predominantly targets establishing security verification frameworks for SoC designs, employing methods such as information flow tracking [3], assertion-based security property verification [4], fuzzing [5], runtime verification monitor [6], and Concolic testing [7]. However, these methods fall short of scalability for large designs. They are not adaptable enough to use in other designs without notable modifications, have limited coverage, and demand significant computational resources [8]. These approaches also necessitate extensive manual analysis of designs, identifying threat models and vulnerabilities to formulate effective security policies for thorough verification. Furthermore, the commercial tools available for security verification often require in-depth knowledge to operate effectively.

To address such limitations of contemporary verification methods, there is a pressing need for a solution that possesses a nuanced understanding of complex SoC designs and can transfer knowledge of one design to another to uncover vulnerabilities effectively but requires the least manual intervention. Large language models (LLMs), with their advanced pattern recognition, natural language understanding, advanced reasoning, and knowledge transfer capabilities, stand out as promising candidates for this task. LLMs can analyze design documentation and codebases to learn the specific characteristics of one design and apply this knowledge to another, potentially revealing security weaknesses that are not immediately apparent. Moreover, LLMs can potentially automate parts of the verification process, significantly reducing the need for manual analysis. However, existing LLMs maintain a fixed context length<sup>1</sup>, which is not long enough to accommodate

<sup>1</sup>Context Length: The number of tokens an LLM can process at once.

a large design. Hence, current LLM-based hardware security solutions focus only on small hardware designs [9] because of token<sup>2</sup> limitation and limited memorization<sup>3</sup> issues.

To put a solution together, we introduce SoCureLLM, an LLM-based hardware security verification technique that solves the abovementioned problems by integrating the methodical partitioning of complex designs, enriching the analysis with contextual summarization for continuity, and employing refined prompting strategies for targeted exploration. Along with the potential to detect hardware vulnerabilities, SoCureLLM also creates a comprehensive security policy database, which is essential for ensuring the security robustness of hardware designs.

The key contributions of our work are summarized as follows:

- SoCureLLM is the first LLM-based hardware security verification framework to handle large-scale SoC designs.
- The framework exhibits scalability and flexibility, effectively addressing the token limitation and memorization problems typically associated with LLMs.
- We introduce an automated vulnerability detection method within SoCureLLM, which thoroughly evaluates designs for potential security flaws.
- SoCureLLM integrates a procedure for creating an extensive security policy database through threat model-based assessments of the large designs.
- The performance of our framework is demonstrated through a rigorous evaluation of various buggy SoC designs, and we compare the performance with existing security verification methods.

In the remainder of this paper, Section II narrates the preliminaries and related works. Later, Section III describes the proposed methodology. Afterward, the experimental setup, a case study, and the results are discussed in Section IV. Finally, Section V concludes the paper.

## II. BACKGROUND & PRIOR WORKS

In Section II-A, we initially provide an overview of the foundational concepts related to LLMs and SoC security to familiarize the reader with the basics. This foundational knowledge is essential for comprehending the proposed methodology and discussions presented in this paper. Subsequently, Section II-B explores the relevant research work.

### A. Preliminaries

#### 1) Large Language Model (LLM)

LLMs, a subset of generative artificial intelligence (GAI), specialize in processing and generating natural language. They are trained on extensive text datasets, which enables them to perform a wide range of linguistic tasks, such as extracting knowledge from the provided information. The transformer architecture [10] serves as the foundation for the

most advanced LLMs such as GPT-3 [11], and GPT-4 [12]. This architecture allows GPT models to generate coherent and contextually relevant text by building upon each token it produces, iteratively refining the output as it progresses through the sequence.

#### 2) In-context Learning & Prompting

In-context learning (ICL) [11] leverages the innate capability of LLMs to adapt their outputs based on given examples or instructions within the input prompt. This allows the models to perform tasks without explicit training or fine-tuning for each new task - a process typically requiring substantial computational resources and time.

Prompting, in turn, is the mechanism by which users communicate with LLMs, providing them with the context or examples they need to learn in context. The quality of the prompt directly influences the success of in-context learning and the performance of the model. A well-engineered prompt can succinctly convey the task requirements and guide the LLM toward the desired output. With this importance in mind, various prompting techniques, such as chain of thought (CoT) [13], have emerged. These techniques aim to optimize how LLMs are guided during inference to produce the most accurate and relevant responses possible.

#### 3) Parameters in LLM models

One significant benefit of LLMs compared to traditional deep learning methods is their ability to control the output generated using various parameters [14]. An example of such a parameter is the temperature, which adjusts the variability in the predictions made by the model. When generating text, an LLM predicts the next word based on a probability distribution that reflects the likelihood of different words appearing in context, as learned from its training data. The temperature setting influences this distribution: A higher temperature results in a broader, more uniform spread, leading to outputs that are more diverse and unpredictable. Conversely, a lower temperature narrows the distribution, making the output more consistent and predictable.

#### 4) SoC Security Vulnerabilities

Hardware security vulnerabilities are weak points in the design of a system that attackers can exploit. For example, during interrupt handling, if the processor is not running at the highest privilege level, any attacker can exploit the interrupt service routine (ISR), potentially gaining control over the entire system. The vulnerabilities might arise unintentionally due to design mistakes or intentionally through rogue employees or compromised 3PIPs, often revealed during the transition from models to physical implementations. CAD tools may inadvertently contribute to these bugs during synthesis as they have not been designed with security in mind. Test and debug infrastructures could also expose the system to post-manufacture security breaches. Common vulnerabilities and exposures (CVE) and common weakness enumeration (CWE) by MITRE are the primary references to classify these issues, offering a comprehensive list of prevalent hardware and software vulnerabilities. Table I shows a portion of the vulnerabilities considered in this paper.

<sup>2</sup>Token: The smallest unit of text LLM processes.

<sup>3</sup>Lack of Memorization: Limitation of LLMs in retaining and recalling information over an extended period.

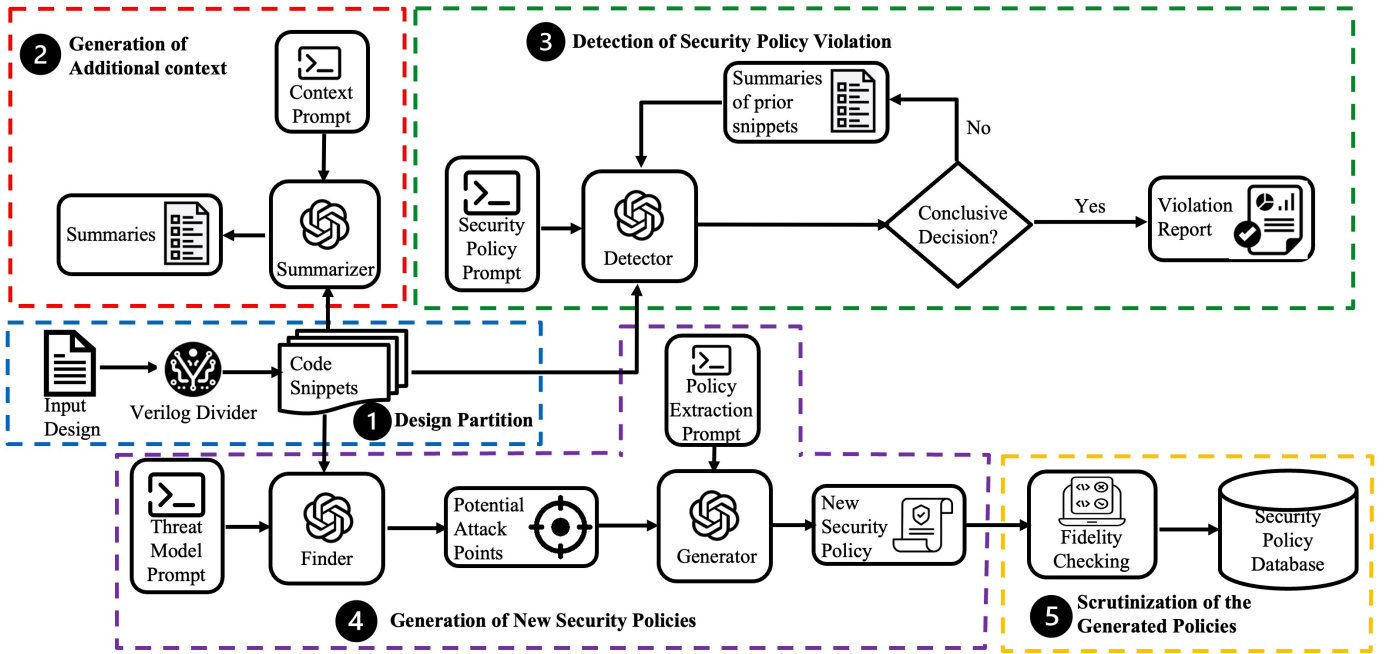


Fig. 1: Overall flow of the proposed SoCureLLM framework

### 5) Security Policy

Modern SoC designs contain sensitive information that requires protection from unauthorized access, guided by principles of confidentiality, integrity, and availability [15]. Security policies translate these principles into practical design constraints, helping IP designers and design integrators in establishing and executing protective measures. Security policies vary, with some ensuring access control, protocol verification, and memory defense, while others guard against information leakage, fault injection, and denial of service. However, formulating these policies for modern SoCs is complex and makeshift, demanding considerable manual labor based on specific customer needs and design details. To address the aforementioned privilege escalation issue, a verification engineer must understand the ISR’s specific implementation in the design under test, as well as the potential consequences of privilege escalation. Crafting such a security policy will be time-intensive due to these detailed requirements.

### B. Prior Works

As the security verification techniques in SoC designs evolve, H. Witharana *et al.* [4] introduced an automated framework for generating security assertions tailored to specific vulnerabilities. However, the effort is limited to a specific set of vulnerabilities. B. Ahmad *et al.* [16] developed security-specific scanners to identify and analyze hardware CWEs and evaluated their performance in different open-source designs. This method still involves manual intervention. Machine learning-based hardware verification methods [17] face challenges such as design dependency, data scarcity, scalability, and efficiency problems, limiting their broader applicability and effectiveness. Recent research efforts [9], [18], [19], [20] using LLM have begun to address the detection and mitigation

of hardware vulnerabilities. D. Saha *et al.* [18] explored the potential capabilities of LLM in different SoC security tasks and mentioned the challenge of LLM-guided methods in handling large hardware design due to the limitation of context length. The study in [21] used LLMs to identify and map security vulnerabilities in SoC designs to relevant CWEs, generate assertions, and enforce security policies, with demonstrated efficacy in open-source SoC benchmarks. W. Fu *et al.* [20] curated a dataset of hardware design defects and remediation steps from open-source projects. They trained medium-sized LLMs to identify hardware bugs but excluded lengthy files exceeding the models’ context limits, leaving large hardware designs unanalyzed. In another work, D. Saha *et al.* [19], identified 16 different security vulnerabilities and weaknesses in small FSM designs using hand-crafted prompts. However, none of these works addressed the challenge of locating security bugs in large designs, which is the main objective of this paper.

## III. PROPOSED SOCURELLM FRAMEWORK

Identifying vulnerabilities early in the design process is paramount to upholding the integrity, confidentiality, and availability of designs. A comprehensive database of security policies greatly enhances the validation of any security verification framework, including the detection of vulnerabilities. Our proposed framework intends to harness the capabilities of LLMs to fulfill these objectives. Within this framework, the LLM adopts different names depending on its function: it is the “Summarizer LLM” when condensing information, the “Detector LLM” when scanning for security policy violations, the “Finder LLM” when locating potential attack points, and the “Generator LLM” when crafting security policies for IPs. The important steps of this framework are

TABLE I: A PORTION OF THE HARDWARE VULNERABILITIES USED IN THIS WORK

CWE ID	Vulnerability Type	Vulnerability Description	CIA Violation	Required Security Policy	Security Implication
CWE-1198	Privilege level	Improper privilege assignment	Confidentiality	The privilege level must revert to the highest level upon returning from a debug session	Unauthorized access control
CWE-269	Privilege level	Improper interrupt handling	Confidentiality	The process of interrupt handling should only occur in the highest privilege level of a processor core	Unauthorized access control
CWE-269	Privilege level	Illegal changes in enable signals	Confidentiality, integrity	Privilege levels should remain unchanged during instruction execution, and updates to read/write enable for control and data registers should not occur during privilege transfers	Unauthorized access control
CWE-250	Privilege level	Illegal virtual page access request	Integrity	The Memory Management Unit (MMU) must consistently uphold the appropriate privilege levels for accessing virtual memory pages	Unauthorized access control
CWE-250	Privilege level	Illegal instruction execution	Confidentiality, integrity	Instructions with the highest privilege level of a process should maintain their privilege status throughout the execution	Unauthorized access control
CWE-1260	Memory access	Memory range overlapping	Confidentiality	The memory address ranges for all security-critical IPs within an SoC need to be distinct and clearly defined	Unauthorized memory access
CWE-284	Memory access	Illegal DMA access	Confidentiality	The DMA controller must be monitored while accessing the protected memory region	Unauthorized memory access
CWE-1245	FSM	Vulnerable FSM encoding	Integrity	When transitioning between two consecutive unprotected states, the Hamming distance between them must be 1	Access to protected states
CWE-1271	Reset related	Important register values were inaccurately cleared during reset	Confidentiality	Upon reset, all critical registers must undergo initialization	Information leakage
CWE-506	Hardware Trojan	A Trojan in the CSR module, causing the Supervisor User Memory Access (SUM) bit of the MSTATUS register to be set to 1, enabling access to user-level virtual pages from supervisor mode	Confidentiality	When supervisor mode does not have access to user-level pages, the SUM bit must be set to 0	Information leakage
CWE-506	Hardware Trojan	Trojan implanted in decoder module for tracking hardware exceptions, causing CPU to halt at threshold	Availability, integrity	The processor needs to monitor unauthorized hardware exceptions	Denial of service (DoS)
CWE-310	Crypto IP	A Trojan within the AES engine, resulting in the encryption process being halted for an uncommon plaintext byte, which acts as a trigger.	Availability, integrity	The AES engine should signal DONE after completing the 10th round of operations for a 128-bit AES cipher key	Denial of service (DoS)
CWE-310	Crypto IP	A Trojan leading to the leakage of the encryption key when encountering a rare plaintext byte	Confidentiality	The AES encryption key must not be exposed	Information leakage
CWE-1244	Debug Module	Illegal JTAG access	Confidentiality	Every debugging session should demand a password, with the password verification system correctly implemented	Information leakage
CWE-1244	Debug Module	Illegal JTAG access	Confidentiality	For every, debug request, there should be a bitwise check after every reset	Information leakage

presented in Figure 1. The following sections give detailed explanations of each step.

#### Step ①: Partitioning of Designs into Smaller Code Snippets:

In this framework, we employ large-scale, open-source buggy SoC designs as our initial input. Since LLM cannot handle such a large volume of tokens simultaneously, our framework adopts a design partitioning approach. A Python-based Verilog code divider splits the design into a series of smaller code snippets based on various rules. Initially, the individual IP modules are separated from each other. Subsequently, each IP is further subdivided into smaller code segments. The code snippets are saved as text files in a local repository. Although the framework allows random segmentation of the design without requiring complete functionality, we have defined a set of rules based on various Verilog constructs to guide the segmentation process. These rules consider factors such as the number of lines of code, the count of ‘always’ blocks, the presence of case statements, and any IP instantiations within the module. For example, consider a control & status register (CSR) IP module comprising 1500 lines of code; it can be divided into 15 distinct code snippets. The first snippet contains all inputs, outputs, and parameters. The second and third snippets correspondingly encapsulate the CSR Read logic and the CSR Write & update logic. The rest of the code is divided into code snippets similarly, depending on the other rules.

#### Step ②: Generation of Additional Context through Summarization:

LLMs often struggle with context retention in interactions. Apart from creating smaller code snippets, we need to preserve and provide the information from earlier segments during the complete analysis. To capture the most significant information from each code snippet, we generate a concise technical natural language summary for every snippet. This stage employs the “Summarizer LLM” to craft these summaries. A carefully constructed context prompt helps the LLM in producing a thorough summary that includes the signal, security, and functional information of the code snippets. These summaries are subsequently stored as text files within a local storage repository. Notably, these summaries are rigorously constrained to a concise fixed length. These condensed summaries play an important role in step 3 as additional context in the framework.

The following textbox presents an example summary derived from a segment of code pertaining to the CSR module. This summary contains essential functional elements, such as procedural blocks, as well as security-related details, including privilege levels. Consider a scenario where the security policy prompt requires verification of a particular policy: after a debug session, the processor must revert to the highest privilege level. When the code snippet is fed to an LLM along with this summary (of a different snippet), it will work as an additional context regarding the processor’s privilege status.

#### Example of summary

This Verilog code snippet consists of three always\_comb blocks: privilege\_check, exception\_ctrl, and wfi\_ctrl. The privilege\_check block checks for privilege violations when performing operations on the Control Status Register (CSR). For a violation, it sets the privilege\_violation flag. The exception\_ctrl block handles exceptions, setting the cause to ILLEGAL\_INSTR and marking the exception as valid if there’s access exception or privilege violation. The wfi\_ctrl block controls the Wait for Interrupt (WFI) instruction. It sets the wfi\_d flag based on the status of mip\_q, debug\_req\_i, irq\_i[1], debug\_mode\_q, csr\_op\_i, and ex\_i.valid. Machine (M) mode is the highest privilege level for the CSR module.

Step ③: Detection of Security Policy Violation: Following the generation of summaries and code snippets, in this phase, our framework proceeds to employ the “Detector LLM”. In this particular step, we employ a predefined policy prompt that contains the security policies to be checked, with the aim of identifying a spectrum of hardware vulnerabilities within the code snippets. To ensure the correctness of the decisions made by the LLM, we follow a specific prompting strategy while crafting the security policy prompt. The details of this prompting strategy are comprehensively described in Section III-A.

In this recursive step, the Detector LLM analyzes each code snippet against a security policy prompt. The LLM then produces a response that may confirm or deny violations or seek more context for a definitive ruling on policy compliance. The results are marked as conclusive if a clear affirmation or negation of the policy violation is confirmed. Otherwise, the decisions are inconclusive. For conclusive findings, the LLM generates an in-depth violation report detailing the findings and explanations.

In contrast, when the decision remains inconclusive, we need to provide additional context from the prior code snippets. Hence, the current code snippet, along with the policy prompt and the summary of the preceding code snippet, is resubmitted to the LLM. This process repeats until the LLM gives a definite answer on policy violations or until it has reviewed all prior summaries. If still undecided, the LLM creates a report stating it is inconclusive and moves on to the next code snippet. For example, if our design is divided into 100 code snippets and we are examining the fourth snippet using the LLM with the policy prompt, it will produce a clear violation report or, if the decision is inconclusive, we provide the LLM with the fourth snippet, policy prompt, and a summary of the third snippet for further analysis. This step is repeated, incorporating summaries of the 2nd and 1st snippets, if needed, until the LLM reaches a conclusive decision. If it remains unresolved after all summaries are reviewed, the LLM compiles an inconclusive report for the fourth snippet before moving on to the fifth.

Given the variability in outcomes across different code snippets - where some may comply with security policies and others may not - the framework adopts a cautious approach. The framework prioritizes critical findings; therefore, if even

---

**Algorithm 1** Detection of Security Policy Violation

---

**Require:**  $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$ : Set of code snippets,  $V$ : Set of vulnerabilities,  $P_v$ : Security policy prompt for  $v \in V$ ,  $P_c$ : Context prompt,  $\mathcal{L}_{\text{summary}}$ : Summarization LLM,  $\mathcal{L}_{\text{detection}}$ : Decision-making LLM,  $\mathcal{C}$ : Conclusive decisions,  $O = \{o_{1,v}, o_{2,v}, \dots, o_{n,v}\}$ : Intermediate detection outcomes,  $R_v$ : Vulnerability report.

**Ensure:**  $\{R_v \mid v \in V\}$ : Set of reports for each vulnerability.

- 1:  $\mathcal{S} \leftarrow \mathcal{D}$
- 2: Store  $\mathcal{S}$  in memory
- 3: **for all**  $v \in V$  **do**
- 4:   **for all**  $s_i \in \mathcal{S}$  **do**
- 5:      $r_i \leftarrow \mathcal{L}_{\text{summary}}(P_c, s_i)$    /\* LLM creates summary \*/
- 6:      $o_{i,v} \leftarrow \mathcal{L}_{\text{detection}}(s_i, P_v)$    /\* LLM performs security check \*/
- 7:     **if**  $o_{i,v} \notin \mathcal{C}$  **then**
- 8:       **for**  $j = i - 1$  **downto** 1 **do**
- 9:          $o_{i,v} \leftarrow \mathcal{L}_{\text{detection}}(s_i, P_v, r_j)$    /\* Includes summary \*/
- 10:        **if**  $o_{i,v} \in \mathcal{C}$  **then**
- 11:         **break**
- 12:    $R_v \leftarrow \max_{\text{severity}}(O)$  /\*Reports as buggy, if detected in any snippet\*/
- 13: **return**  $\{R_v \mid v \in V\}$

---

**Algorithm 2** Generation of Security Policy

---

**Require:**  $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$ : Set of code snippets,  $T$ : Set of threat models,  $\mathcal{L}_{\text{finder}}$ : Finder LLM,  $\mathcal{L}_{\text{generator}}$ : Generator LLM,  $P_t$ : Threat model prompt,  $P_e$ : Policy extraction prompt,  $\mathcal{F}_c$ : Fidelity Checking.

**Ensure:**  $SP = \{sp_{1,1}, sp_{1,2}, \dots, sp_{n,m}\}$ : Set of security policies approved after fidelity check.

- 1: Initialize  $SP = \emptyset$
- 2: **for all**  $s_i \in \mathcal{S}$  **do**
- 3:   Initialize  $AP = \emptyset$
- 4:   **for all**  $t \in T$  **do**
- 5:      $AP_i \leftarrow \mathcal{L}_{\text{finder}}(s_i, P_t)$  /\* LLM identifies potential attack points\*/
- 6:      $AP \leftarrow AP \cup AP_i$    /\* Aggregate attack points \*/
- 7:   **for all**  $ap \in AP$  **do**
- 8:      $sp_{i,ap} \leftarrow \mathcal{L}_{\text{generator}}(s_i, ap, P_e)$    /\* LLM creates security policy \*/
- 9:      $passed \leftarrow \mathcal{F}_c(sp_{i,ap})$    /\* Manual scrutiny evaluates the policy \*/
- 10:    **if**  $passed$  **then**
- 11:      $SP \leftarrow SP \cup \{sp_{i,ap}\}$    /\* Add only approved policies to SP \*/
- 12: **return**  $SP$

---

a snippet exhibits noncompliance or a security flaw, the whole design is marked as vulnerable. In the context of our previous example, even if 1 out of 100 code snippets is found to have issues after a thorough iterative assessment, the framework would still categorize the overall design as potentially vulnerable. The algorithm 1 concisely depicts the hardware vulnerability detection flow (step ①, ② and ③) of the proposed framework.

**Step ④: Generation of New Security Policies:** In this stage, we require the inclusion of code snippets as input to the ‘Finder LLM’. Distinct from the policy prompt approach, we employ an alternative strategy. We craft a threat model

prompt containing a curated selection of threat models capable of investigating breaches in integrity, confidentiality, or availability within the design. This prompt is meticulously designed to encompass various scenarios stemming from the chosen threat models. The threat models in consideration include the following categories: information leakage, denial of service, confidentiality attack, privilege escalation, access control violation, and unauthorized memory access. Upon submitting these inputs, the Finder LLM formulates a response that identifies potential points of vulnerability within the design. These vulnerabilities can manifest as sensitive signals, specific conditions, or distinct case statements. The ‘Generator LLM’ is supplied with potential attack points in conjunction with the policy extraction prompt. This policy extraction prompt directs the Generator LLM to transform these potential attack points into actionable security policies.

**Step ⑤: Scrutinization of Generated Policies:** The generated policies are then subjected to a thorough examination. It is possible that certain policies derived may coincide with the security policies mentioned in the security policy prompt. Additionally, some of the policies might be impractical to implement. Therefore, this process of scrutinization filters out such policies, retaining only the viable ones in an extensive security policy database. Algorithm 2 concisely represents the security policy generation aspect (step ①, ④ and ⑤) of the proposed framework.

### A. Prompting Strategies

As discussed in Section II-A2, since the performance of an LLM is significantly influenced by the prompting technique used, our framework is designed with customized prompting strategies to enhance performance.

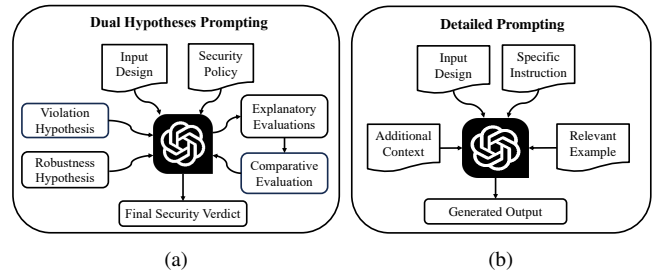


Fig. 2: Different prompting strategies used in SoCureLLM: Overview of dual hypotheses prompting (a) and detailed prompting (b)

While formulating the security policy prompt in step ③ of the proposed framework, we adopt a novel prompting strategy named the ‘dual hypotheses prompting approach’, shown in Figure 2a. In this prompting strategy, the LLM is tasked first to consider two assumptions: violation and robustness hypotheses. In the ‘violation hypothesis’, the LLM is asked to assume that there is a violation of the given security policy within the design and generate plausible explanations for such violation. Then, in the ‘robustness hypothesis’, the LLM is prompted to consider the contrary, focusing on the absence of



vulnerabilities by highlighting robust security features of the design. By contrasting these two scenarios, the strategy facilitates a comprehensive evaluation, comparing the explanations to ascertain which is more convincing. The strategy concludes by evaluating the most logical explanation to determine the presence or absence of the violation of the security policy. This approach essentially mirrors scientific hypothesis testing, adapted for the domain of SoC security. It enhances critical evaluation by contrasting scenarios of risk, leading to more informed and rational decision-making in SoC security.

Furthermore, we follow another prompting technique (shown in Figure 2b) in designing the context, threat model, and policy extraction prompts. In this method, the LLM is guided by specific instruction, additional context, relevant examples, and self-scrutiny instruction along with the input design to produce enriched output. In this case, specific instructions ensure the process adheres to defined criteria. Additional context and relevant examples enrich decision-making by providing depth and situational awareness, while self-scrutiny instruction is a critical introspective step in which the LLM reviews and optimizes its own processes to ensure accuracy and relevance.

#### IV. EXPERIMENTS

In order to confirm the efficacy of the proposed framework SoCureLLM, an extensive experiment has been performed on multiple SoC designs. Section IV-A presents the experimental setup. Next, Section IV-B explains the methodology through a case study. Finally, Section IV-C describes the experimental results with analysis.

##### A. Experimental Setup

Although any LLM can be used in the SoCureLLM framework, we applied the GPT-4 [12] API within our framework as the core LLM. We focused on detecting security policy violations in three vulnerable RISC-V SoCs: hack@dac 2018 (PULPissimo core [22]), TrustHub benchmark, and HOST 2022 (Ariane core). The hack@dac 2018 SoC is a renowned SoC with known bugs, often used as a benchmark in various verification methods. Since the GPT-4 model was trained with data up to September 2021, it might be familiar with the hack@dac 2018 SoC. To circumvent issues related to the availability of design information, we have opted for two other SoCs, HOST 2022 and the TrustHub benchmark, as they are mainly closed-source. We utilized three additional open-source RISC-V SoCs for the security policy generation flow: CVA6, lowRISC(ibex), and CV32E40P. A Python-based tool was developed to segment Verilog code for analysis. Simulations were conducted using Verilator and Cadence Incisive, while Cadence Jaspergold SPV and FPV tools were employed for formal verification.

##### B. Case Study

To evaluate the effectiveness of our proposed framework for detecting vulnerabilities in SoCs, we selected three RISC-V-based SoCs as detailed in the preceding section. This case

#### Dual Hypothesis Prompting (C.1)

**//Input Code Snippet & Summaries will be added to the prompt automatically//**

*Security Policy: The debug module's password-checking mechanism must be logically correct, and all related signals should be cleared during the reset state.*

Consider two scenarios: First, assume that the security policy is breached in the provided code snippet and the technical summaries. Write an analysis based on this assumption.

Second, assume that the security policy is adhered to in the provided code snippet and the technical summaries. Write an analysis based on this assumption.

Compare the two analyses and identify the more logical one. If you reach to a conclusive decision, provide a detailed explanation based on your analysis. If your decision is inconclusive, you will request additional summaries from the previous snippets sequentially.

Your response should strictly maintain the following structure:

**Analysis under Breach of Security Policy Assumption:**

# add your response here

**Analysis under Adherence to Security Policy Assumption:**

# add your response here

**Final Verdict:**

# add your response here

# "The decision is conclusive" or "The decision is inconclusive"

study illustrates the process of identifying a vulnerability within the debug module of the TrustHub benchmark and the hack@dac 2018 SoC. The bug is located within the JTAG TAP controller module and is characterized by a flawed password verification mechanism. This mechanism is crucial for any debugger utilizing the JTAG module, with the presumption that an attacker could mimic a debugger to bypass this password verification. The vulnerability lies in the "correct" counter, which tallies the number of correct password bits during each attempt. The counter increases when the input password bit aligns with the actual password bit, unlocking the module when the counter matches or exceeds the total number of password bits. However, during the reset state, this "correct" counter was not cleared to zero due to the bug, leaving the module perpetually unlocked after a single successful entry, thereby creating a potential backdoor for attackers.

SoCureLLM initially segments the entire SoC into multiple code snippets based on the integrated IPs and peripherals. Subsequently, these IPs are broken down into smaller code snippets, which are saved as local text files alongside a technical summary for each. The process is not repeated

### Response from GPT-4 without correct Summary (C.2)

#### Analysis under Breach of Security Policy Assumption:

Retention of previous state information (`passchk`, `correct`, `bitmask`) in other non-reset states can be exploited if these variables are not appropriately re-initialized when expected. An attacker could manipulate the transition conditions by ensuring that `correct` hits the specific value, thereby setting `passchk` to 1 without a legitimate authentication sequence.

#### Analysis under Adherence to Security Policy Assumption:

If the system enters the reset state, it should explicitly clear signals like `passchk`, `correct`, `bitmask`, and others to adhere to the policy that all related signals should be reset. This ensures that any remaining or leftover signal values from prior operations that could jeopardize security are eliminated each time the machine is reset.

#### Final Verdict:

A definitive conclusion on vulnerability cannot be reached without further details on other parts of the system. Additional snippets or a broader view of the full implementation would be necessary to determine compliance or breach of the stated security policy conclusively.

**The decision is inconclusive.**

for the same SoC for other security policies. To identify security policy violations (Step 3 in our framework), each snippet is analyzed by the detector LLM with a security policy prompt, employing the dual hypothesis prompting technique. The specific prompt used in the proposed framework to detect the vulnerability of the JTAG module mentioned above is shown in Textbox C.1. This prompt instructs the detector LLM to evaluate based on two hypotheses: detecting security policy violations and ensuring compliance with the security policy. From these assessments, the detector LLM determines whether it can make a definitive decision. If the decision is conclusive, it produces a detailed report of the violation. If the outcome is inconclusive, the SoCureLLM framework supplies the summaries of previous code snippets sequentially until a definitive conclusion is achieved. The prompt additionally ensures that the detector LLM consistently delivers structured responses. In the case of this vulnerability, one snippet encompasses the actual implementation of the password-checking logic, and a separate snippet contains the reset state that includes the bug. When the snippet with the password-checking logic is analyzed without any summary of the previous code snippets, the response from the detector LLM, as displayed in Textbox C.2, does not yield a definitive conclusion. Nonetheless, it outlines the potential implications of both violating and adhering to the security policy as outlined in the prompt. However when the detector LLM reviews the technical summary describing the reset state, where the bug resides, along with the provided code snippet, it conclusively identifies the snippet as vulnerable, as indicated in Textbox

### Response from GPT-4 with correct Summary (C.3)

#### Analysis under Breach of Security Policy Assumption:

Because `correct` is not cleared, an attacker who successfully unlocks the JTAG module once can bypass the password-checking mechanism in subsequent attempts. The counter value remains above the conditional check (`correct == 32'h8000_FFFF`), automatically setting `passchk` to 1, which should only happen after a successful password check.

This persistent `correct` signal effectively disables the password-checking mechanism after one successful entry, making the module vulnerable to unauthorized access without further password verification. The module stays unlocked indefinitely, a severe security risk allowing persistent, unauthorized access.

#### Analysis under Adherence to Security Policy Assumption:

With the additional context provided, it becomes evident that despite resetting `passchk`, failing to reset `correct` is a significant oversight. This non-compliance with the security policy undermines any other security measures implemented in the code.

The initial analysis assumed the other signals might be correctly managed based on the snippet. However, knowing now that `correct` is not reset, it is clear that the implementation does not adhere to the essential security policy of clearing all sensitive signals during a reset.

#### Final Verdict:

Based on the combined analysis of both code snippets and the new summary provided, the code is definitively vulnerable. The failure to reset the `correct` signal allows for an easily exploitable loophole that can leave the module unlocked indefinitely, providing unauthorized access without subsequent password verification. This represents a serious security flaw requiring immediate correction to comply with security policies and ensure the integrity of the module.

**The decision is conclusive.**

C.3. Notably, the security policy violation was so evident that the detector LLM found no grounds to support policy adherence. This case study emphasizes the importance of the technical summaries in pinpointing security policy violations and mitigating the limitations of LLM memorization, affirming the value of feedback and segmentation-based examination within SoC designs through SoCureLLM, thus enhancing the accuracy of security assessments.

### C. Evaluation

We examined three buggy SoCs containing 68 hardware vulnerabilities in the security policy violation evaluation. Due to space limitations, we cannot detail each vulnerability, encouraging readers to go through [22], [23] for comprehensive information about the vulnerabilities and the security policies used in the policy prompt. A general overview of the bugs can be found in Section II-A4.

As discussed in Section III, for the identification of security bugs, our proposed method incorporates several key



TABLE II: Performance comparison between open-ended security assessment and proposed SoCureLLM in the detection of security bugs

Design	IP	LoC	Open Ended Security Assessment				Security Policy Violation Assessment (SoCureLLM)				
			# bugs in the module	# bugs detected	# bugs detected successfully	TPR	TNR	# bugs detected	# bugs detected successfully	TPR	TNR
Hack@Dac2018	Debug Unit	715	7	5	4	0.571	0.950	8	6	0.857	0.900
	GPIO	408	5	3	2	0.400	0.955	4	4	0.800	1.000
	CSR	1510	3	1	0	0.000	0.958	2	2	0.667	1.000
	RISC-V Core	14635	7	3	2	0.286	0.950	9	5	0.713	0.800
	AXI Interface	810	1	0	0	0.000	1.000	1	1	1.000	1.000
	Crypto Modules	11606	4	2	1	0.250	0.957	2	2	0.500	1.000
	<b>Total</b>	<b>29684</b>	<b>27</b>	<b>14</b>	<b>9</b>	<b>0.333</b>	<b>0.963</b>	<b>26</b>	<b>20</b>	<b>0.741</b>	<b>0.956</b>
Trust-Hub Benchmark	CSR	1510	10	4	3	0.300	0.944	8	8	0.800	1.000
	RISC-V Core	14635	4	5	2	0.500	0.870	5	4	1.000	0.957
	Decoder	1418	4	2	1	0.25	0.957	4	3	0.750	0.957
	MMU	519	2	0	0	0.000	1.000	1	1	0.500	1.000
	PMP	278	1	0	0	0.000	1.000	1	1	1.000	1.000
	AES	12624	4	2	2	0.500	1.000	3	3	0.750	1.000
	AXI Interface	810	2	1	1	0.500	1.000	2	1	0.500	0.960
		<b>Total</b>	<b>31794</b>	<b>27</b>	<b>14</b>	<b>9</b>	<b>0.333</b>	<b>0.970</b>	<b>24</b>	<b>21</b>	<b>0.778</b>
HOST 2022	CSR	1510	4	3	2	0.500	0.900	3	3	0.750	1.000
	Crypto Module	8133	5	2	2	0.400	1.000	5	4	0.800	0.890
	RISC-V Core	14635	4	3	2	0.500	0.900	5	3	0.750	0.800
	memory unit	235	1	0	0	0.000	1.000	1	1	1.000	1.000
		<b>Total</b>	<b>24513</b>	<b>14</b>	<b>8</b>	<b>6</b>	<b>0.423</b>	<b>0.952</b>	<b>14</b>	<b>11</b>	<b>0.786</b>

steps, including partitioning the design, summarizing, and implementing novel prompting techniques. To demonstrate the effectiveness of these enhancements, we compared our results with an open-ended, LLM-driven, zero-shot bug identification approach that identifies bugs without specific context about the design in question. In this open-ended approach, we supplied the LLM with a buggy design and posed a broad question aimed at identifying general bugs. Table II highlights the effectiveness of our framework in identifying vulnerabilities through security policy violations compared to open-ended security assessment. The comparison indicates that SoCureLLM has a true positive rate (TPR) of 76.47%, significantly outperforming the open-ended method’s 35.29%, showcasing its higher efficiency in accurately detecting vulnerabilities present in the design. This highlights the critical role of specific prompting and summarization in identifying security policy violations. The open-ended verification approach lacked defined prompts and summarization, leading to its inability to detect most bugs. In addition, the framework achieves a higher true negative rate (TNR) of 96.46%, showing its reliability in confirming the absence of bugs and minimizing false positives. The lines of code (LoC) in the table demonstrate our framework’s capability to process large-scale designs, effectively addressing the token limitation and memorization issues in LLMs.

Figure 3 compares traditional security verification methods (formal verification, simulation, static code analysis) and SoCureLLM. The results indicate that SoCureLLM outperforms other verification approaches in bug detection, specifically detecting 16.18% more bugs than the closest competitor (simulation). Although simulation or static code analysis can catch some bugs, they are generally slower than the other methods. Moreover, it takes significant manual effort to prepare the test stimuli for the process. Regarding formal verification, it can pinpoint vulnerabilities when there are clear deviations from specified requirements, such as value mismatches or parameter

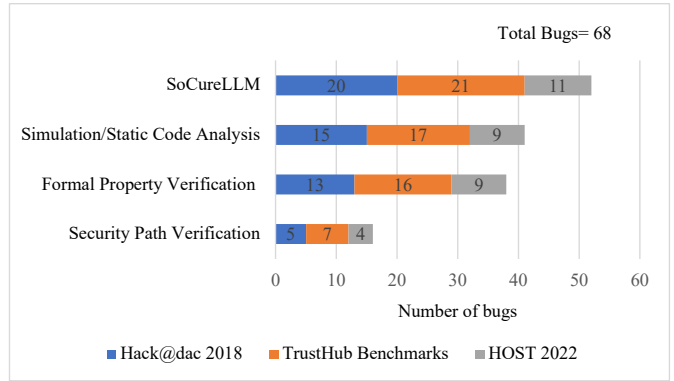


Fig. 3: Comparison of our proposed framework with the contemporary security verification techniques in terms of the number of bugs detected successfully in the SoCs under test

inconsistencies (i.e., memory range overlaps, privilege escalations, and reset-related issues). SoCureLLM is faster and more user-friendly as it only needs design specifications to find vulnerabilities, unlike formal tools that take more time and need detailed knowledge of the specifications and implementation. For instance, autoSVA [24], a formal verification approach, requires approximately one hour to evaluate the MMU module of the Ariane core, while SoCureLLM completes the analysis, including verifying all security policies, in just 15 minutes.

Some vulnerabilities, such as hardware Trojans, are inherently difficult to detect due to their stealthy nature, posing challenges even for SoCureLLM. However, a potential approach might involve fine-tuning a custom LLM with a substantial dataset of hardware Trojan examples, enabling it to effectively identify such threats to some extent. Future work could also focus on inter-modular bugs by integrating an additional mechanism to analyze interactions between modules, enabling the detection of these complex vulnerabilities. In terms of expenses, a complete analysis of the hack@dac 2018 SoC for

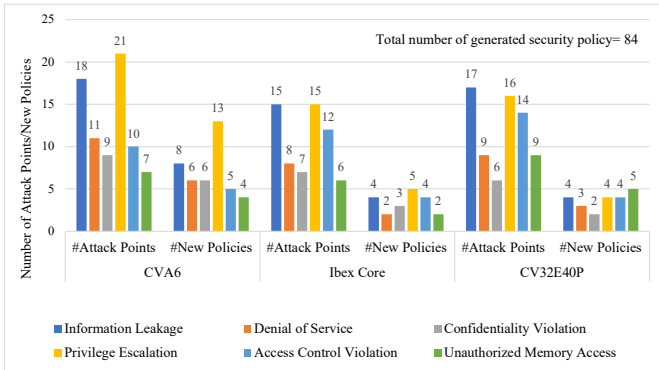


Fig. 4: Number of potential attack points for different threat model analyses in each SoC and newly generated security policies after scrutinization.

vulnerability detection costs just around \$7.5, a figure that is significantly economical compared to the costs of time-to-market and commercial licensing tools. Furthermore, this cost is continually decreasing with the advancement of LLM technology.

Figure 4 illustrates the quantity of potential weak points and the security policies derived from chosen threat models for three SoCs. The weak points are consistently about twice the number of security policies, attributed to multiple attack points leading to similar types of security policies. The SoCureLLM develops security policies by spotting potential vulnerabilities, but manual review is crucial to refine them. This step is necessary to eliminate overly unfeasible ones like those for brute-force attacks and dismiss broad ones such as suggesting the use of physical unclonable functions (PUFs) or cryptomodels. Despite these issues, the manual scrutiny process used here is still more efficient than creating security policies for each threat model from the ground up. A noticeable decrease in the count of security policies across different designs is observed due to the scrutinization process eliminating repetitive policies from each SoC. Eventually, our framework generated 84 new security policies from the designs. A sample list of the generated security policies for the lowRISC (ibex) controller module is shown in the corresponding textbox.

#### Sample of generated security policies

- Monitor the 'debug\_mode\_o' signal to ensure it does not leak through side channels.
- Review state transitions for potential escalation vectors, especially around 'current\_priv\_lvl\_i'.
- Implement rate-limiting and sanity checks on the 'irq\_req\_ctrl\_i' signal to prevent IRQ flooding.
- Implement strict access control checks on debug and control registers.
- Audit the control flow for any unauthorized bypasses or weak checks.
- Implement bounds checking for memory accesses and handle exceptions for 'data\_misaligned\_i'.

#### D. Discussion

When evaluating our method compared to traditional static bug detection techniques, it becomes evident that our approach achieves superior precision scores. We also compared our approach with a recently published state-of-the-art bug identification technique [25]. Our method demonstrated a significant improvement in precision, achieving approximately 0.95 compared to their reported precision of around 0.70. Both results were evaluated using the hack@dac 2018 Pulpissimo buggy SoC. Although static methods are well known for their use in hardware design and verification, they also fail to provide absolute assurance in terms of security bug detection. In contrast, our proposed method not only outperforms the leading security bug identification techniques but also offers distinct advantages that distinguish it from conventional static analysis.

SoCureLLM offers a streamlined and automated implementation process, making it simpler and more efficient than traditional static methods, which often require extensive manual effort. The ease of implementation extends to adaptability as well; our approach is inherently flexible to accommodate new vulnerabilities without requiring significant changes to the underlying framework. This adaptability is due to the consistent prompting strategies and implementation criteria that we have established, which remain the same regardless of the type or nature of the vulnerability being addressed. On the other hand, static bug identification methods often necessitate specific modifications for each new design, which can be time-consuming and resource-intensive. Furthermore, SoCureLLM employs rule-based partitioning and summarization techniques, making it highly effective for analyzing large-scale designs. Hence, SoCureLLM is designed to be compatible with any LLM, whether open-source or closed-source, ensuring its adaptability across different platforms.

A key advantage of our method is its scalability. Through an iterative approach, we decompose the design into manageable components, enabling a thorough examination of each part and efficient handling of complex systems. This scalability ensures our approach can adapt to larger, more intricate designs without sacrificing performance, addressing a significant limitation of many static methods.

A crucial aspect of maintaining the reliability of our method involves controlling the variability of the responses generated by LLM. To this end, we have set the temperature value used in the experiments to be consistently low, as discussed in Section II-A3. This setting minimizes randomness in the outputs, ensuring that similar prompts will produce similar responses. This consistency is vital for achieving more deterministic behavior from the LLM, which is a key requirement for any verification method that aims to be integrated into a formal verification flow.

#### V. CONCLUSION

SoCureLLM offers a flexible, scalable solution to hardware security verification for large-scale SoCs, employing an LLM-based framework to effectively navigate the complexities

of modern SoC designs. It outperforms traditional methods by efficiently detecting security vulnerabilities and enriching security policy databases. SoCureLLM not only solves the challenges in the existing hardware verification techniques but also addresses the token limitation and memorization constraints of traditional LLMs. The future direction of this research involves fine-tuning an LLM specifically for hardware security verification to enhance its precision and effectiveness within the domain.

#### ACKNOWLEDGMENT

This research was funded by the U.S. National Science Foundation (NSF) under the Faculty Early Career Development (CAREER) Program, Grant No. 2339971.

#### REFERENCES

- [1] Z. Kenjar, T. Frassetto, D. Gens, M. Franz, and A.-R. Sadeghi, "VOLTpwn: Attacking x86 processor integrity from software," in *29th USENIX Security Symposium (USENIX Security 20)*, 2020, pp. 1445–1461.
- [2] H. Khattri, N. K. V. Mangipudi, and S. Mandujano, "Hsdl: A security development lifecycle for hardware technologies," in *2012 IEEE International Symposium on Hardware-Oriented Security and Trust*, 2012, pp. 116–121.
- [3] A. Ardeshiricham, W. Hu, J. Marxen, and R. Kastner, "Register transfer level information flow tracking for provably secure hardware design," in *Design, Automation & Test in Europe Conference & Exhibition (DATE), 2017*, 2017, pp. 1691–1696.
- [4] H. Witharana, A. Jayasena, A. Whigham, and P. Mishra, "Automated generation of security assertions for rtl models," *J. Emerg. Technol. Comput. Syst.*, vol. 19, 2023.
- [5] T. Trippel, K. G. Shin, A. Chernyakhovsky, G. Kelly, D. Rizzo, and M. Hicks, "Fuzzing hardware like software," in *USENIX Security Symposium*, 2022, pp. 3237–3254.
- [6] A. Kassem and Y. Falcone, "Detecting fault injection attacks with runtime verification," in *Proc. of the 3rd ACM Workshop on Software Protection*, 2019, p. 65–76.
- [7] X. Meng, S. Kundu, A. K. Kanuparthi, and K. Basu, "Rtl-contest: Concolic testing on rtl for detecting security vulnerabilities," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 41, no. 3, pp. 466–477, 2021.
- [8] S. R. Rajendran, S. Tarek, B. M. Hicks, H. M. Kamali, F. Farahmandi, and M. Tehranipoor, "Hunter: Hardware underneath trigger for exploiting soc-level vulnerabilities," in *2023 Design, Automation Test in Europe Conference Exhibition (DATE)*, 2023, pp. 1–6.
- [9] B. Ahmad, S. Thakur, B. Tan, R. Karri, and H. Pearce, "On hardware security bug code fixes by prompting large language models," *IEEE Transactions on Information Forensics and Security*, vol. 19, pp. 4043–4057, 2024.
- [10] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [11] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 1877–1901.
- [12] OpenAI, "Gpt-4 technical report," 2023. [Online]. Available: <https://arxiv.org/pdf/2303.08774.pdf>
- [13] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Adv. in Neural Information Processing Systems*, vol. 35, pp. 24 824–24 837, 2022.
- [14] [Online]. Available: <https://platform.openai.com/docs/api-reference/chat/create>
- [15] A. Basak, S. Bhunia, and S. Ray, "A flexible architecture for systematic implementation of soc security policies," in *2015 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, 2015, pp. 536–543.
- [16] B. Ahmad, W.-K. Liu, L. Collini, H. Pearce, J. M. Fung, J. Valamehr, M. Bidmeshki, P. Sapiecha, S. Brown, K. Chakrabarty *et al.*, "Don't cweat it: Toward cwe analysis techniques in early stages of hardware design," in *Proc. of the 41st IEEE/ACM Int. Conf. on Computer-Aided Design*, 2022, pp. 1–9.
- [17] Z. Pan and P. Mishra, "A survey on hardware vulnerability analysis using machine learning," *IEEE Access*, vol. 10, pp. 49 508–49 527, 2022.
- [18] D. Saha, S. Tarek, K. Yahyaeei, S. K. Saha, J. Zhou, M. Tehranipoor, and F. Farahmandi, "Llm for soc security: A paradigm shift," *IEEE Access*, vol. 12, pp. 155 498–155 521, 2024.
- [19] D. Saha, K. Yahyaeei, S. Kumar Saha, M. Tehranipoor, and F. Farahmandi, "Empowering hardware security with llm: The development of a vulnerable hardware database," in *2024 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*, 2024, pp. 233–243.
- [20] W. Fu, K. Yang, R. G. Dutta, X. Guo, and G. Qu, "Llm4sechw: Leaving domain-specific large language model for hardware debugging," *Asian Hardware Oriented Security and Trust (AsianHOST)*, 2023.
- [21] S. Paria, A. Dasgupta, and S. Bhunia, (2023) Divas: An llm-based end-to-end framework for soc security analysis and policy-based protection.
- [22] G. Dessouky, D. Gens, P. Haney, G. Persyn, A. Kanuparthi, H. Khattri, J. M. Fung, A.-R. Sadeghi, and J. Rajendran, "HardFails: Insights into Software-Exploitable hardware bugs," in *USENIX Security Symposium*, 2019, pp. 213–230.
- [23] S. Tarek, H. Al Shaikh, S. R. Rajendran, and F. Farahmandi, "Benchmarking of soc-level hardware vulnerabilities: A complete walkthrough," in *2023 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*. Los Alamitos, CA, USA: IEEE Computer Society, jun 2023, pp. 1–6.
- [24] M. Orenes-Vera, A. Manocha, D. Wentzlaff, and M. Martonosi, "Autosva: Democratizing formal verification of rtl module interactions," in *ACM/IEEE Design Automation Conference (DAC)*, 2021, pp. 535–540.
- [25] S. S. Miftah, S. Kundu, A. Mordahl, S. Wei, and K. Basu, "Rtl-spec: Rtl spectrum analysis for security bug localization," in *2024 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 2024, pp. 171–181.